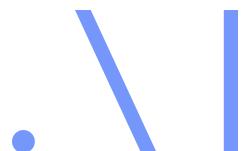


Object Recognition

Georgia Gkioxari



Content

- A Brief Intro
- Visual Recognition: Problem Definitions & Tasks
 - Semantic Segmentation
 - FCN
 - Object Detection
 - Fast(er) R-CNN
 - Instance Segmentation & Pose Prediction
 - Mask R-CNN
- ImageNet in 1H: Important Lessons
- Useful Tips

Content

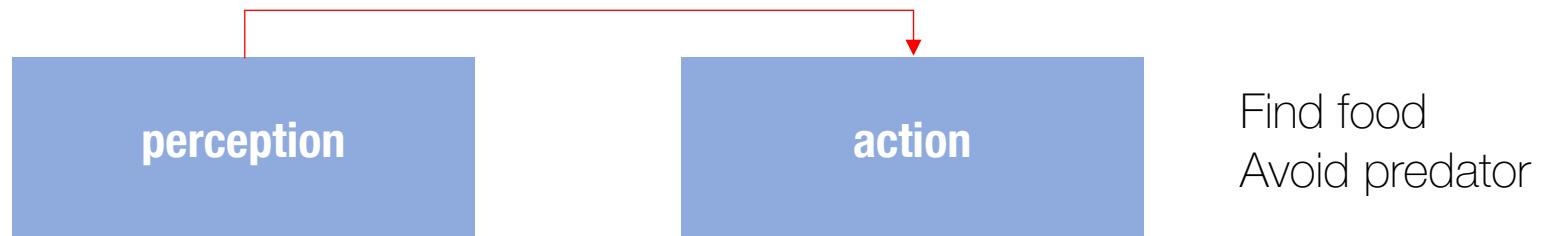
- A Brief Intro
- Visual Recognition: Problem Definitions & Tasks
 - Semantic Segmentation
 - FCN
 - Object Detection
 - Fast(er) R-CNN
 - Instance Segmentation & Pose Prediction
 - Mask R-CNN
- ImageNet in 1H: Important Lessons
- Useful Tips

Why Perception?

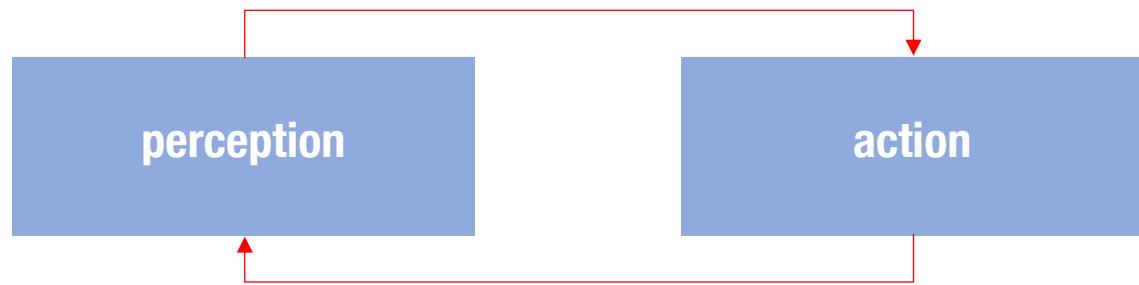
perception

action

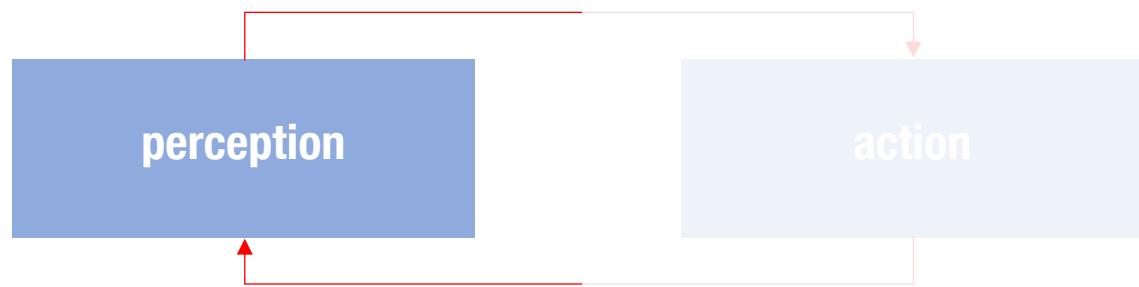
Why Perception?



Why Perception?



Why Perception?



A Brief Intro

A Brief Intro

Invariance: Transformation changes have no effect

$$f(g(x)) = f(x), \forall x$$

Equivariance: Transformation changes in input lead to corresponding changes in output

$$\exists M_g, f(g(x)) = M_g(f(x)), \forall x$$

A Brief Intro

Classification desires **invariant** representations: output a label

Detection/Segmentation desires **equivariant** representations:

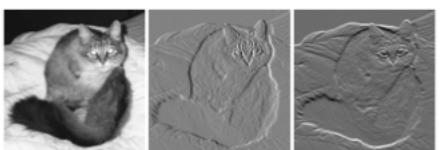
- Translated Object → Translated Mask
- Scaled Object → Scaled Mask

A Brief Intro

- Convolutions are translation-equivariant
- Fully ConvNets (FCN) are translation-equivariant
- ConvNets become translation-invariant due to fully connected layers or global pooling layers

A Brief Intro

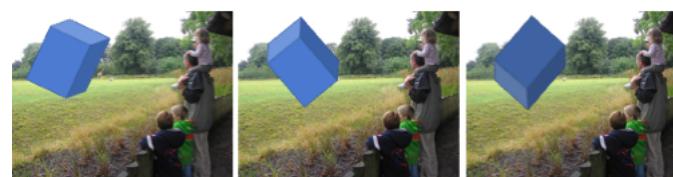
Sources of Variation



Photometric



Translation



Rotation, rigid pose



Scale



Occlusion (self, other)



Deformation (local, articulated)



Intra-class
(appearance, structural)

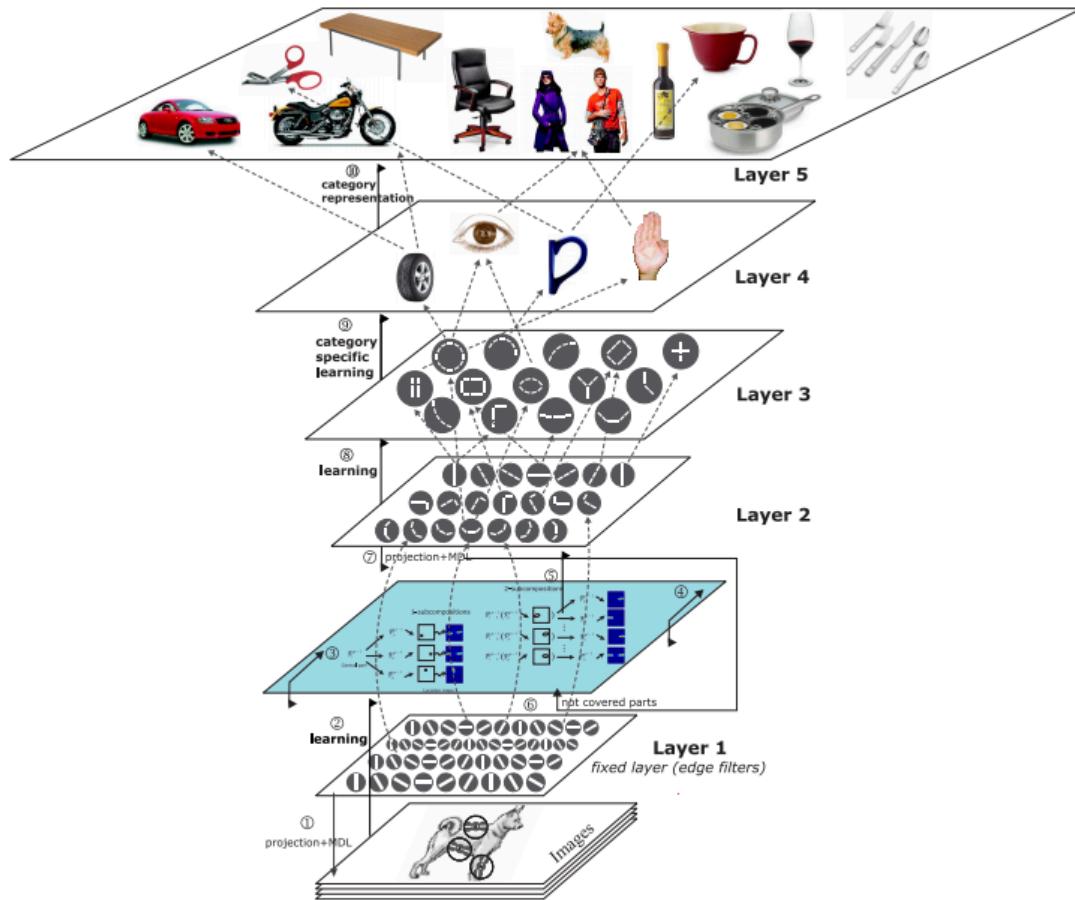


Context
(clutter, spurious corr.)

A Brief Intro

Recognition is the science and engineering of models
that have effective invariance and equivariance properties

A Brief Intro



- Scenes are composed of “stuff” and “things” (objects)
- Objects are composed of parts
- Parts are composed of subparts
- Subparts are built from edges

A Brief Intro

When designing a neural network, think how

- You are modeling equivariance/invariance

A Brief Intro

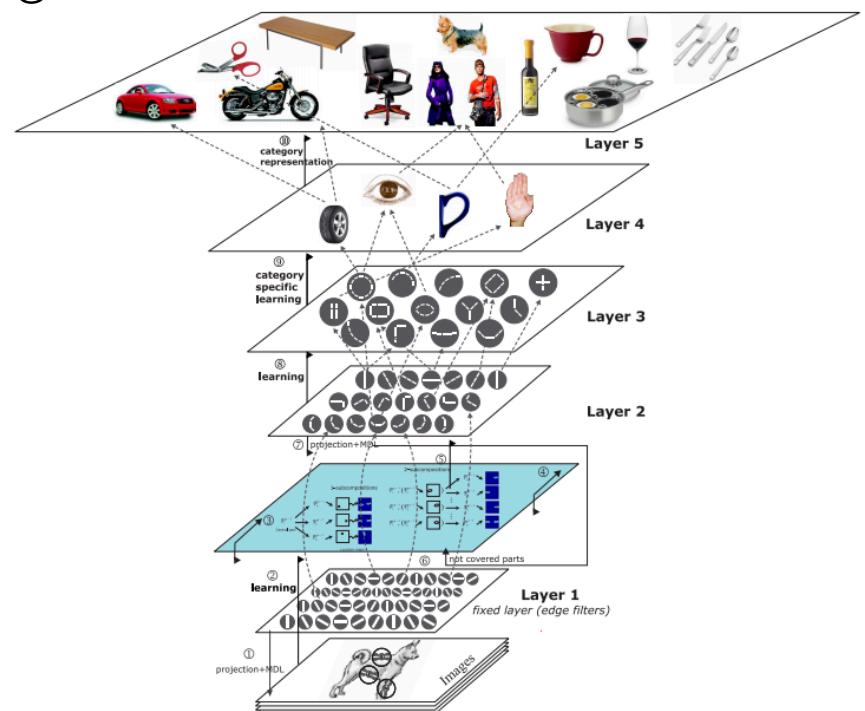
When designing a neural network, think how

- You are modeling equivariance/invariance
- You are leveraging the compositional structure

A Brief Intro

When designing a neural network, think how

- You are modeling equivariance/invariance
- You are leveraging the compositional structure



A Brief Intro

When designing a neural network, think how

- You are modeling equivariance/invariance
- You are leveraging the compositional structure
- You are modeling the computational task as an algorithm

A Brief Intro

When designing a neural network, think how

- You are modeling equivariance/invariance
- You are leveraging the compositional structure
- You are modeling the computational task as an algorithm
- You are organizing the computations so that the solution is efficient

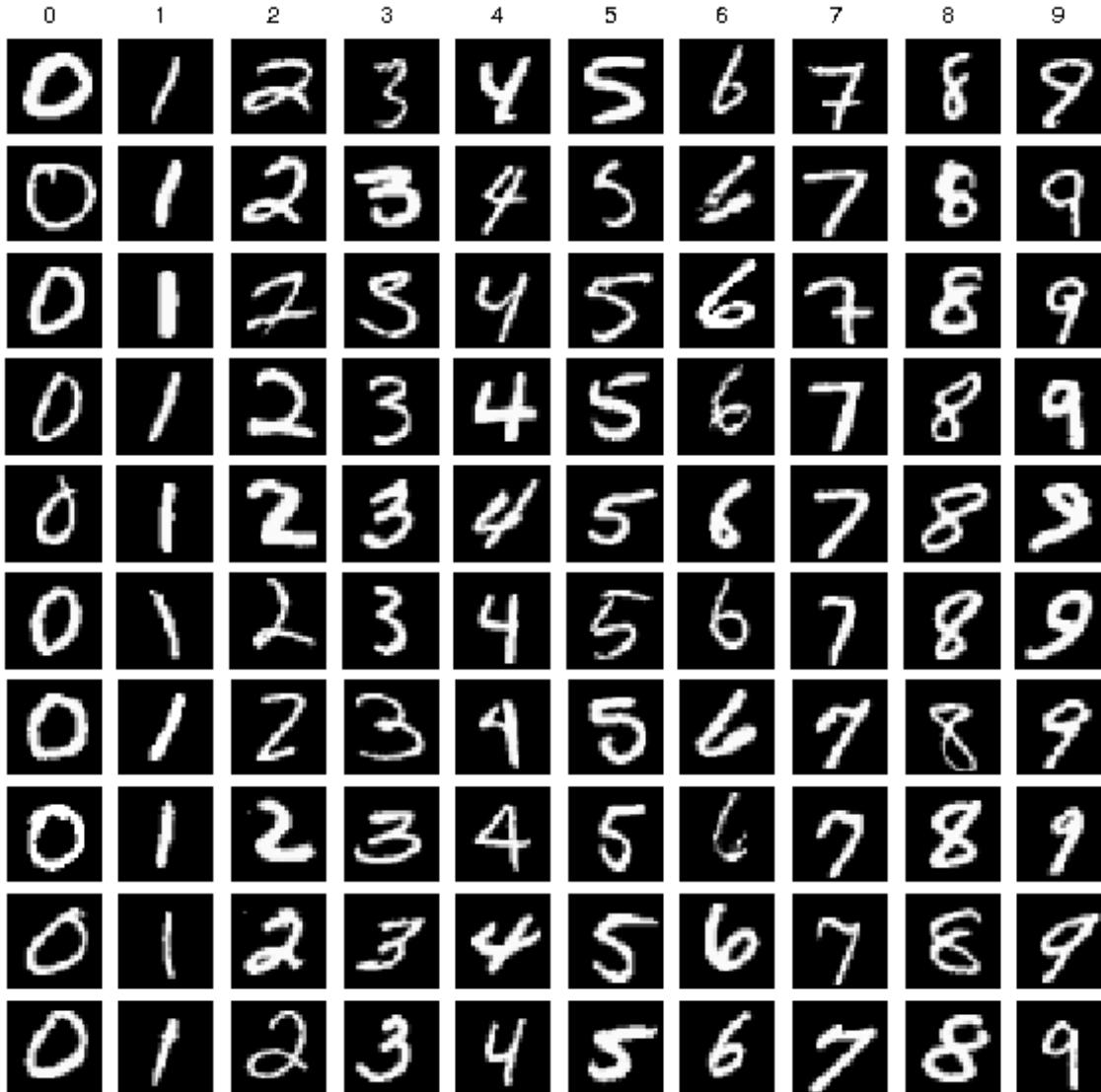
A Brief Intro

When designing a neural network, think how

- You are modeling equivariance/invariance
- You are leveraging the compositional structure
- You are modeling the computational task as an algorithm
- You are organizing the computations so that the solution is efficient

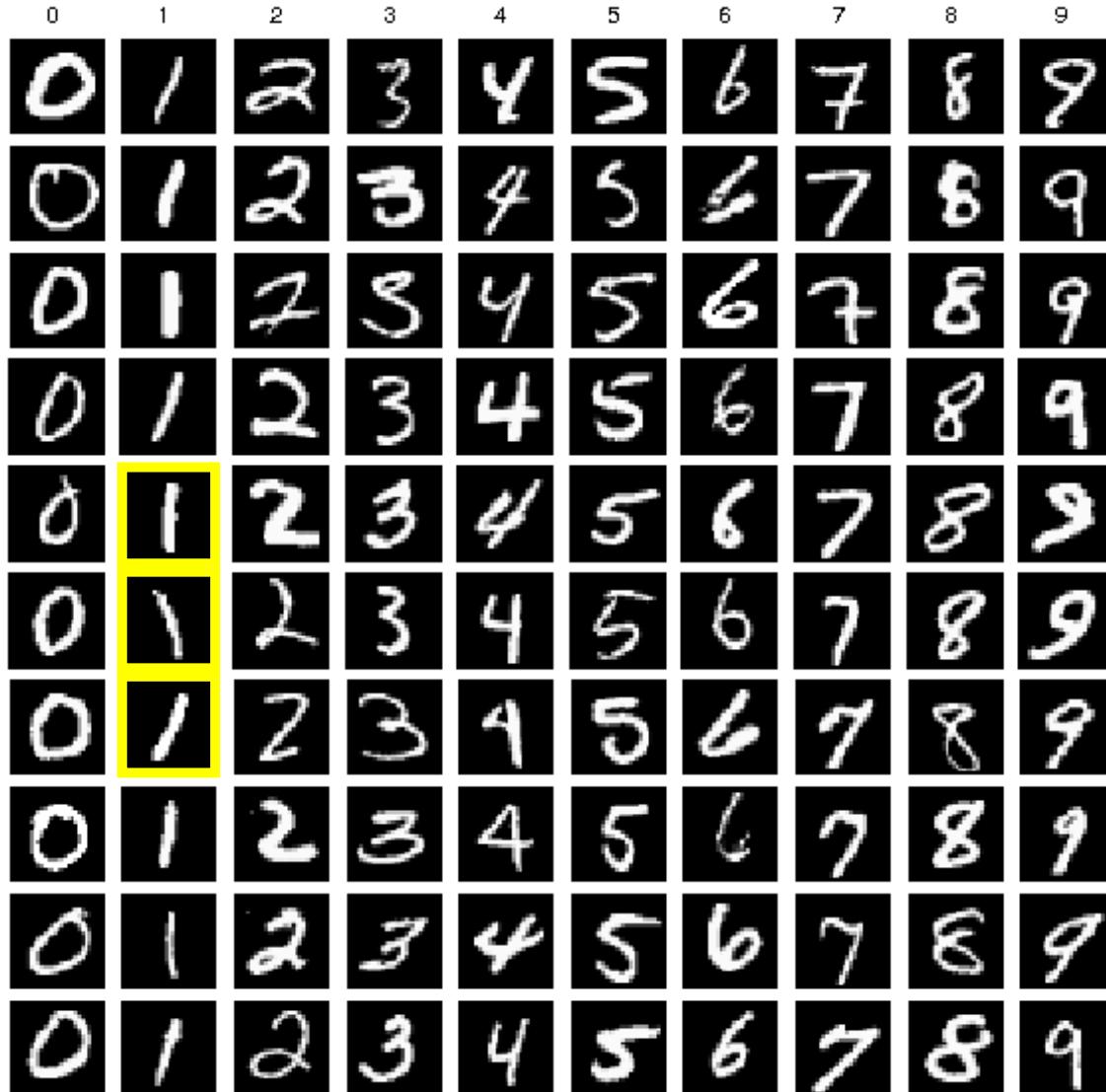
→ **Success**

Isolated Handwritten Digit Recognition



Task: classify into categories $\{0, 1, \dots, 9\}$

Isolated Handwritten Digit Recognition

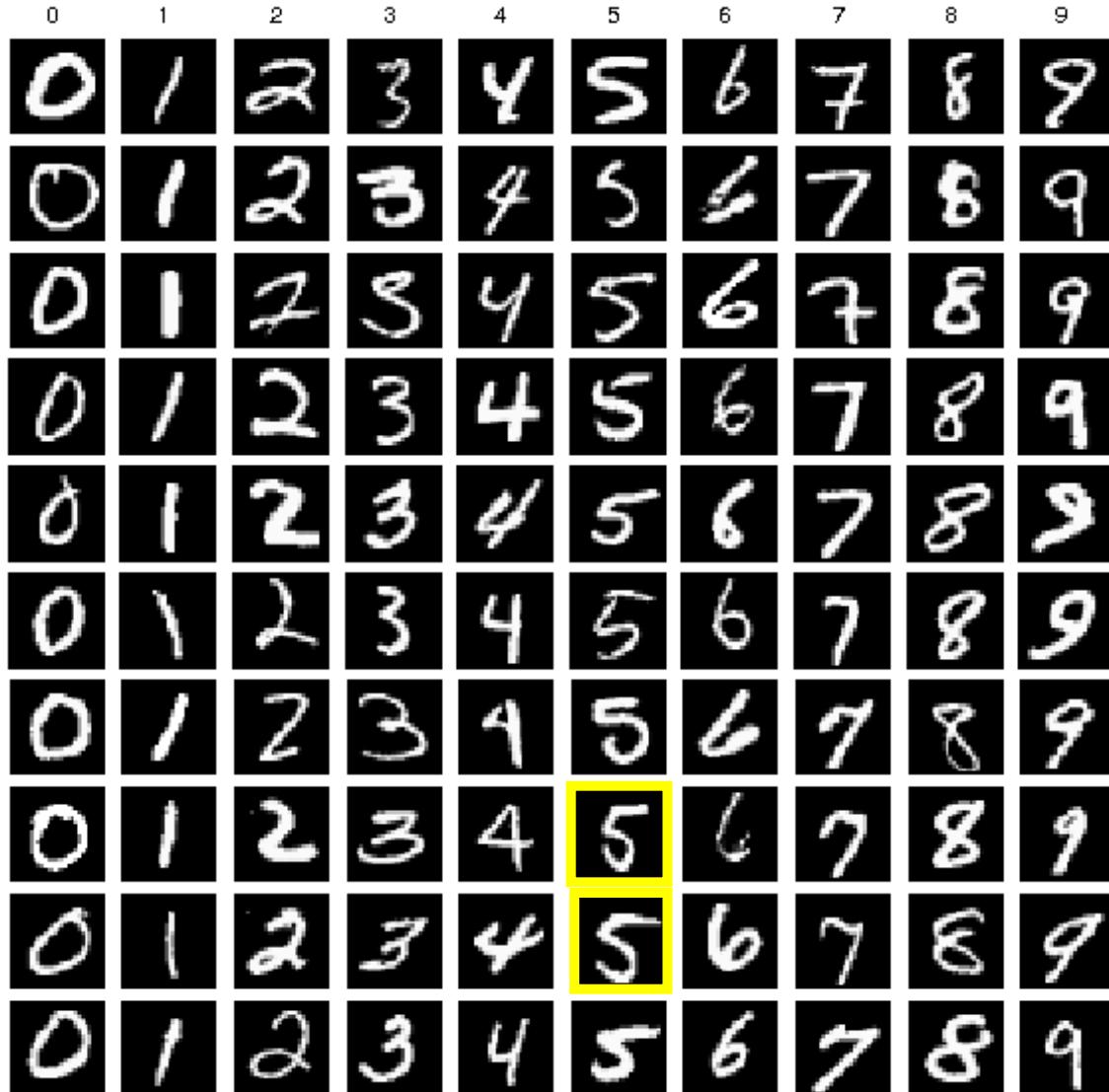


Task: classify into categories $\{0, 1, \dots, 9\}$

Sources of variation:

- Small rotations (in plane)

Isolated Handwritten Digit Recognition

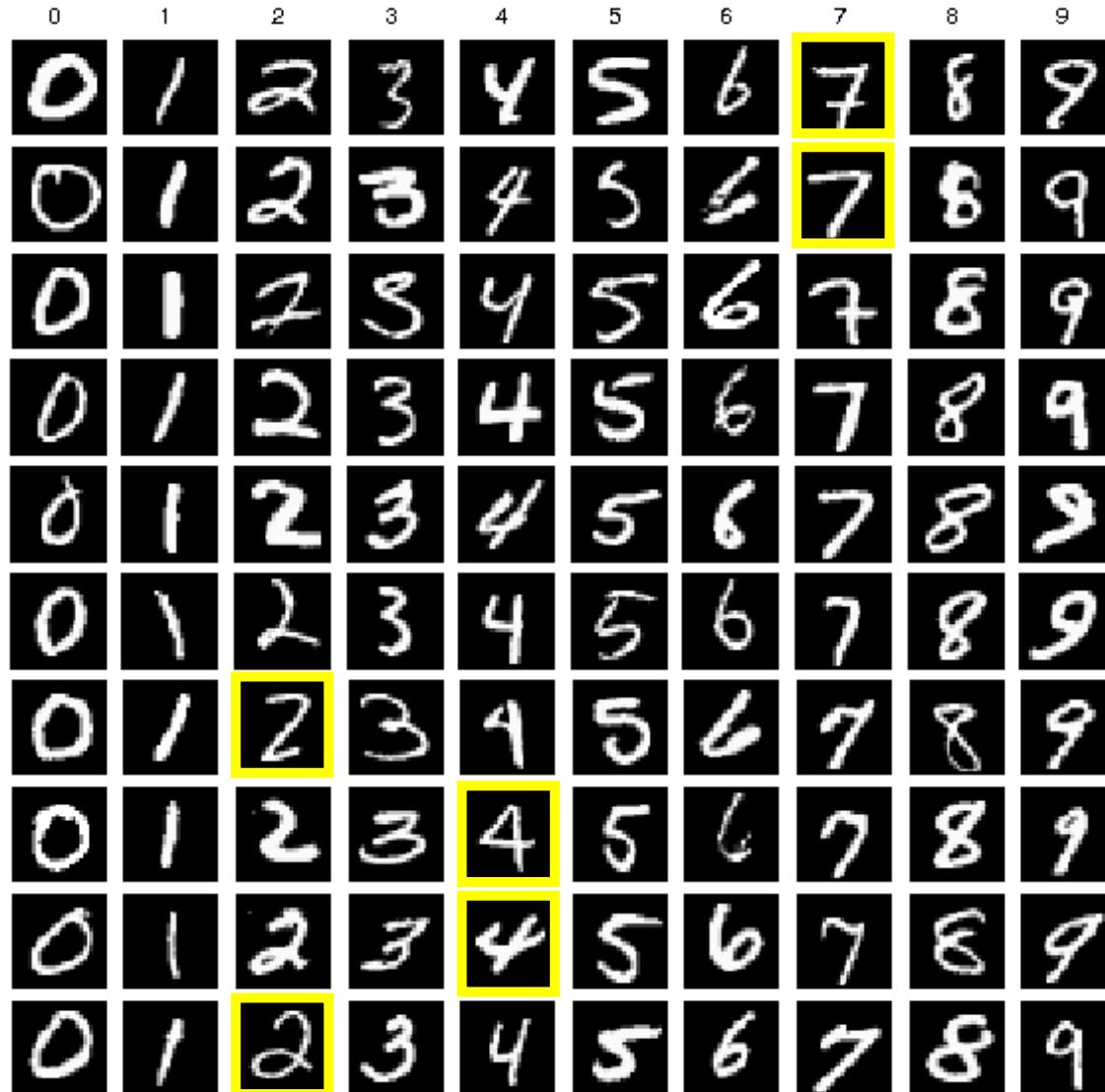


Task: classify into categories $\{0, 1, \dots, 9\}$

Sources of variation:

- Small rotations (in plane)
- Local deformation

Isolated Handwritten Digit Recognition

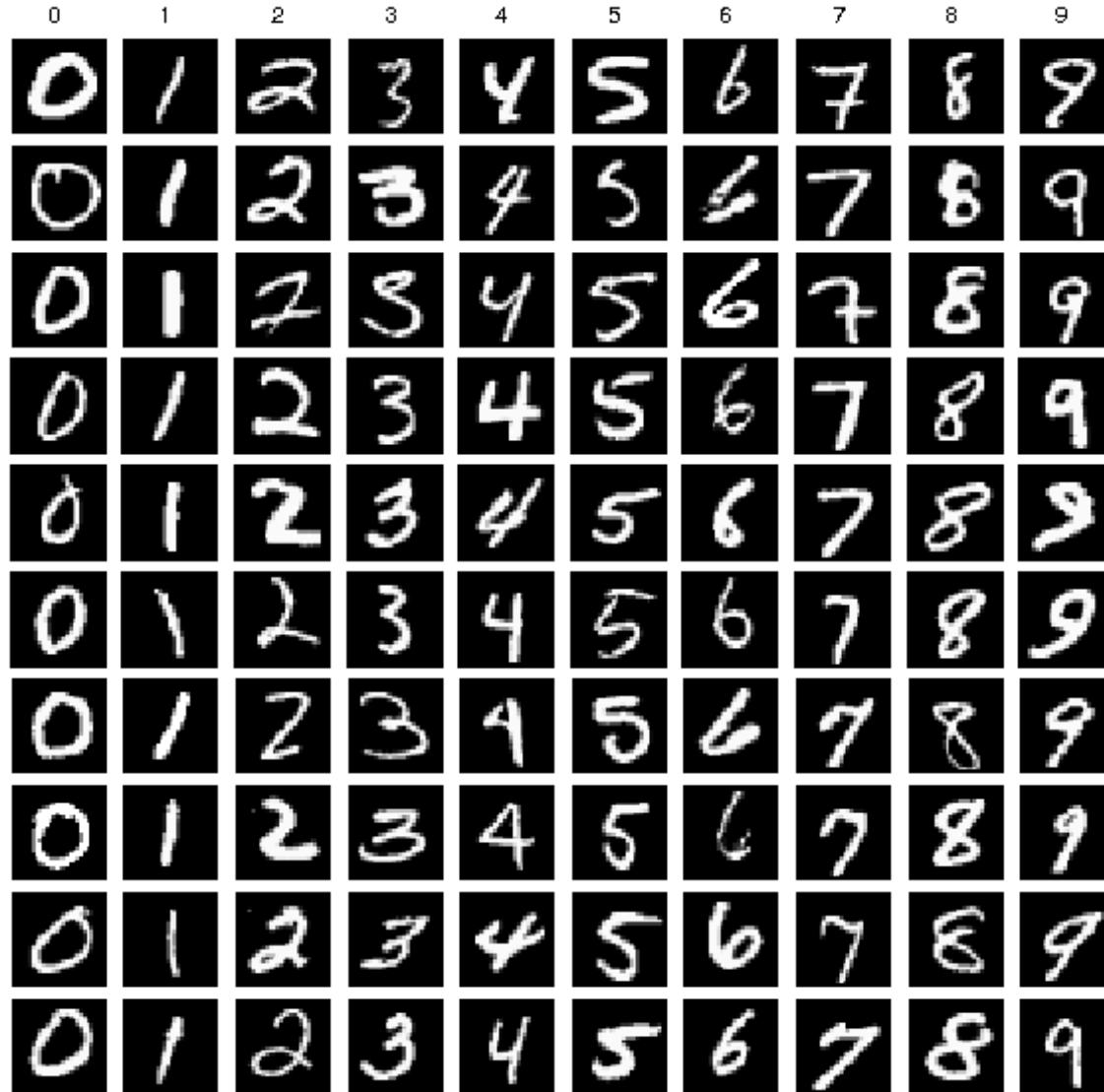


Task: classify into categories $\{0, 1, \dots, 9\}$

Sources of variation:

- Small rotations (in plane)
- Local deformation
- Intra-class variation (style; see 4s and 7s)

Isolated Handwritten Digit Recognition



Task: classify into categories $\{0, 1, \dots, 9\}$

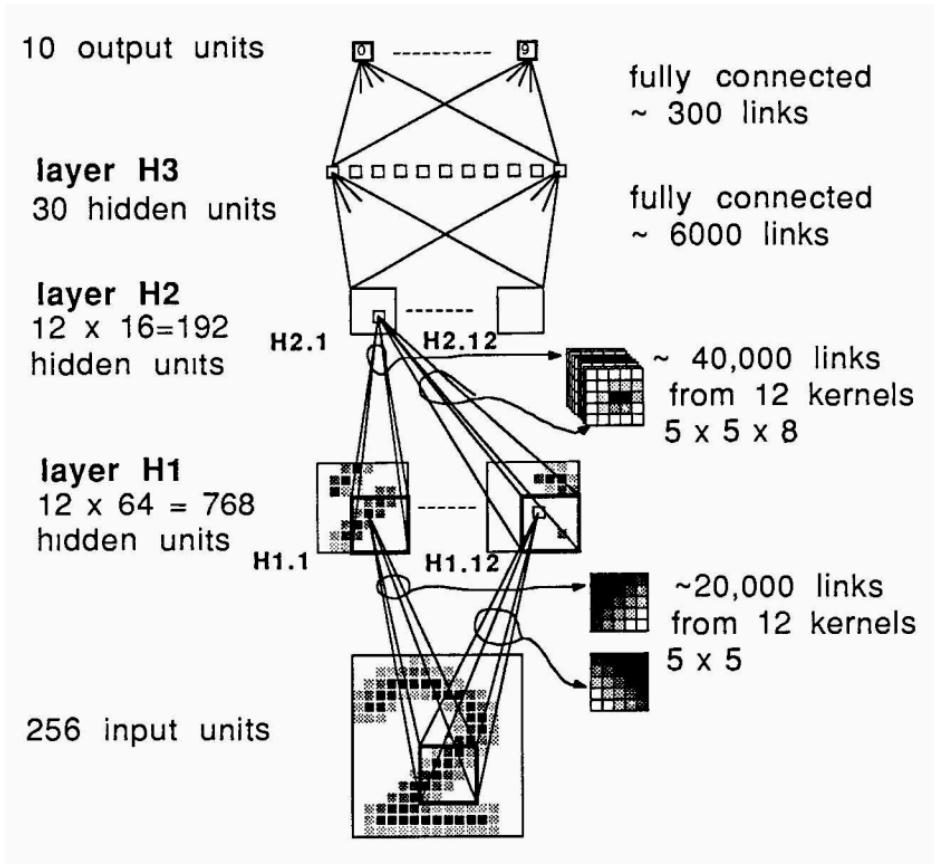
Sources of variation:

- Small rotations (in plane)
- Local deformation
- Intra-class variation (style; see 4s and 7s)

Non-sources of variation:

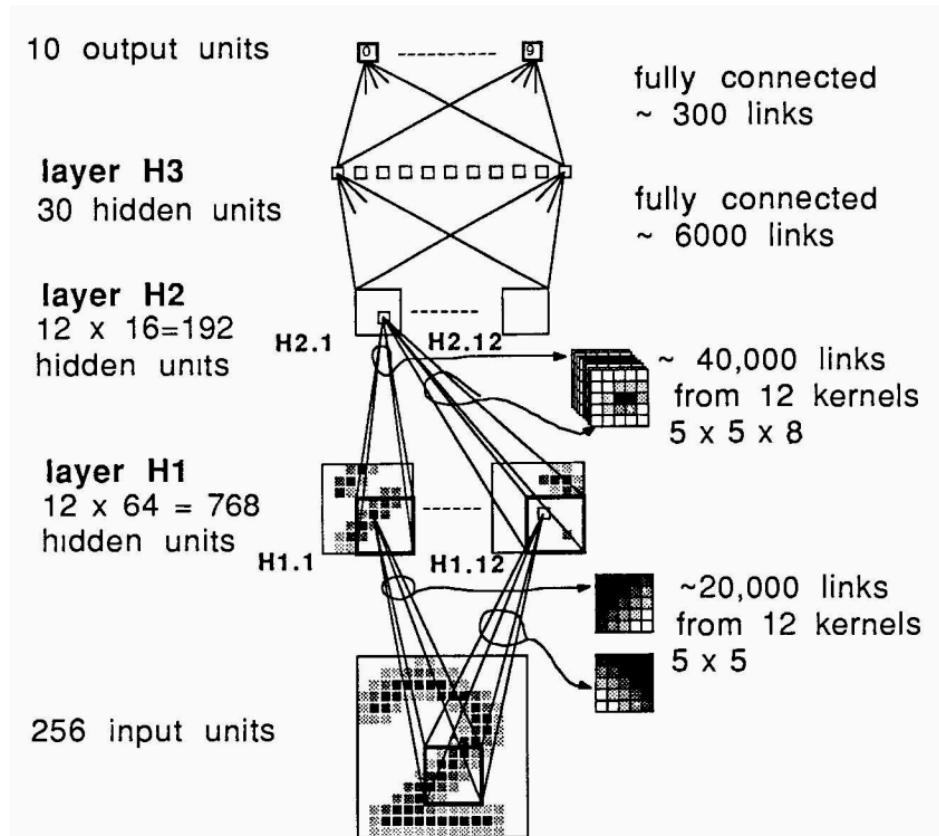
- Photometric
- Translation
- Rotation (out of plane)
- Scale
- Occlusion
- Articulated deformation
- Context

A Brief Intro: Digit Classification



Backpropagation Applied to Handwritten Zip Code Recognition. LeCun et al. 1989.

A Brief Intro: Digit Classification



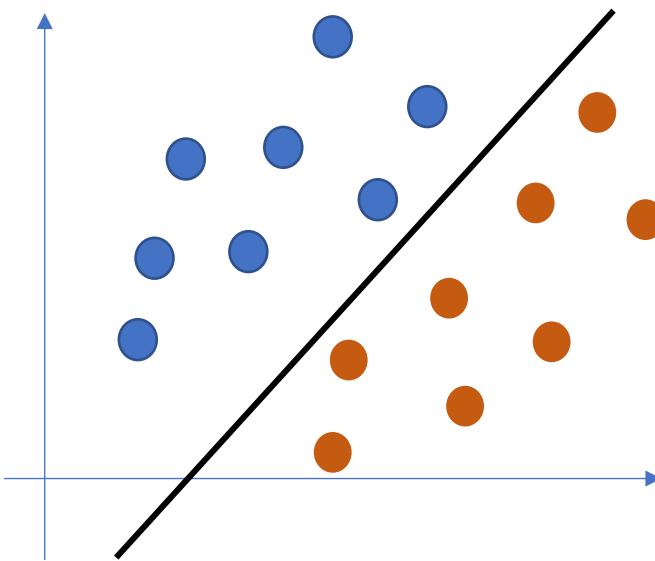
- Two convolutional layers
- One fully connected layer
- Spatial contraction via strided convolution (subsampling)
- Training via SGD with backpropagation

A Brief Intro: Optimization

Assume a dataset $D = \{(x_i, y_i)\}_{i=0}^{N-1}$

We want to learn a function f such that $y_i = f(x_i), \forall i$

E.g. binary classification $y \in \{-1, 1\}$ and $x \in \mathbb{R}^2$



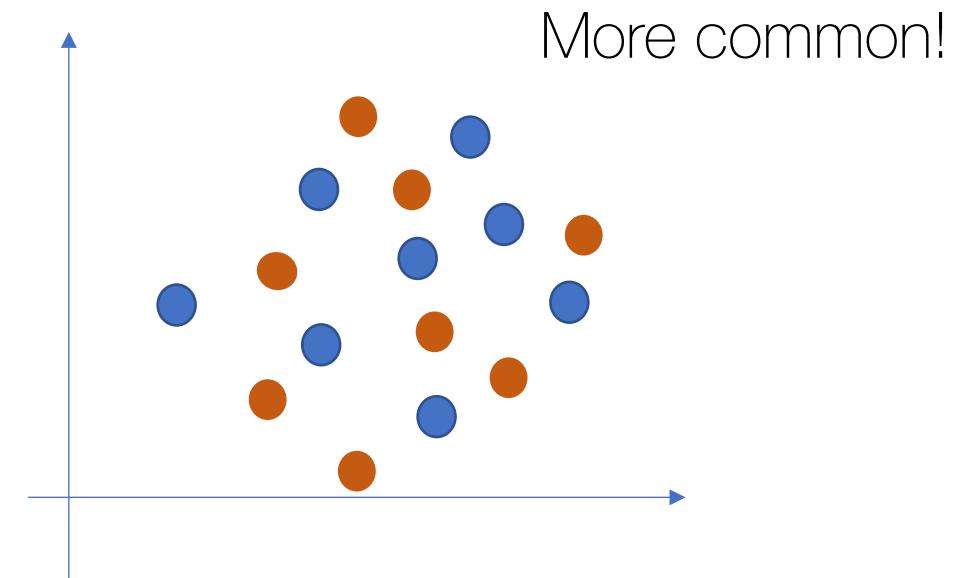
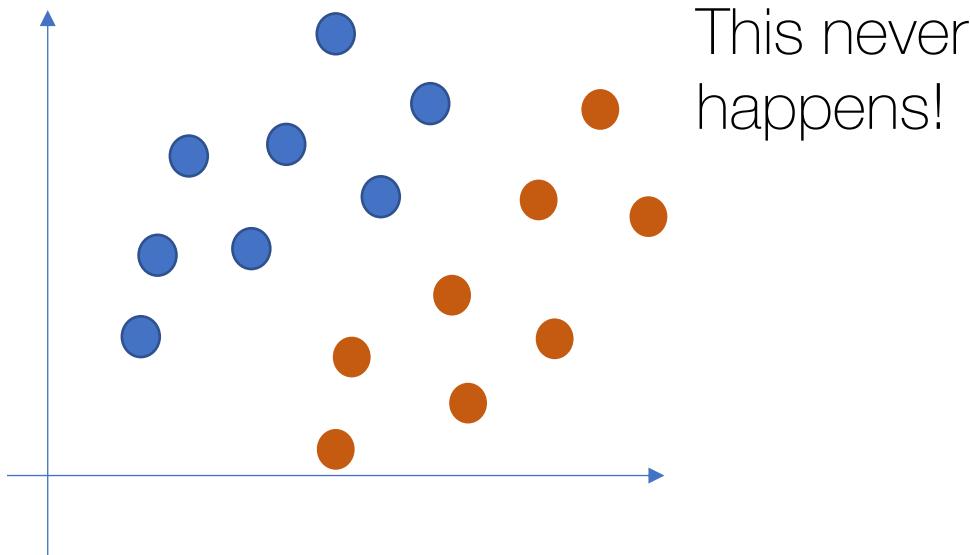
Linearly separable so f is a linear function of x

A Brief Intro: Optimization

Assume a dataset $D = \{(x_i, y_i)\}_{i=0}^{N-1}$

We want to learn a function f such that $y_i = f(x_i), \forall i$

E.g. binary classification $y \in \{-1, 1\}$ and $x \in \mathbb{R}^2$



A Brief Intro: Optimization

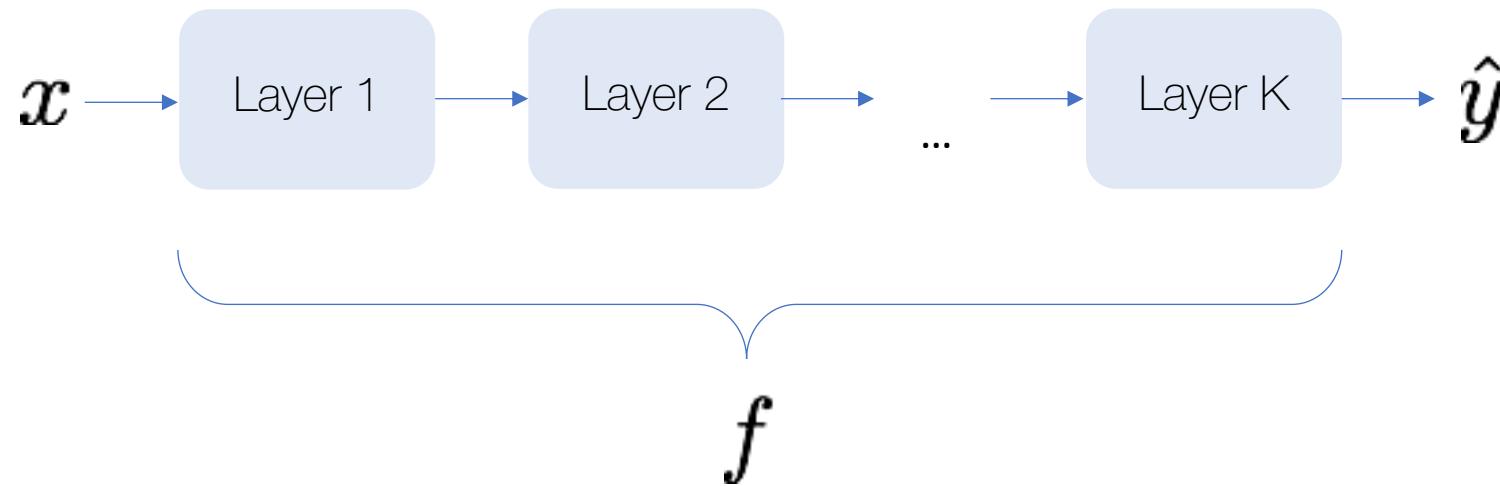
Assume a dataset $D = \{(x_i, y_i)\}_{i=0}^{N-1}$

We want to learn a function f such that $y_i = f(x_i), \forall i$

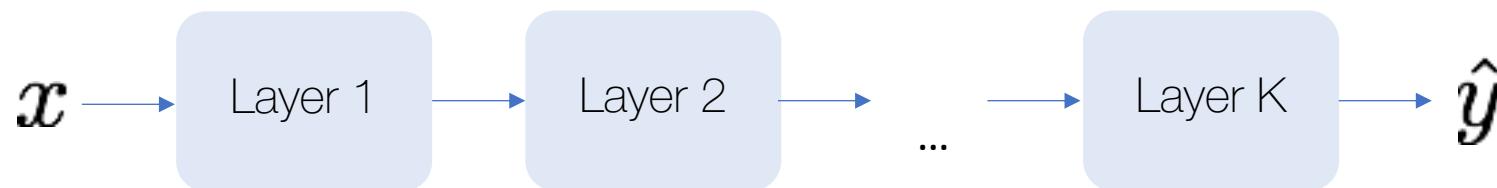
Deep Learning

Find a non-linear function approximator f that best describes our dataset D

A Brief Intro to Deep Learning



A Brief Intro to Deep Learning

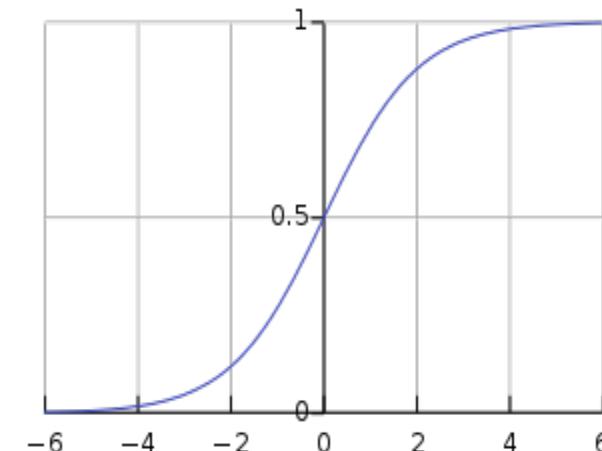


Learnable parameters w_j

Layer j performs a specific operation,
e.g.

$$out_j = \sigma(w_j \cdot in_j)$$

Activation functions σ insert non-linearities,
e.g.

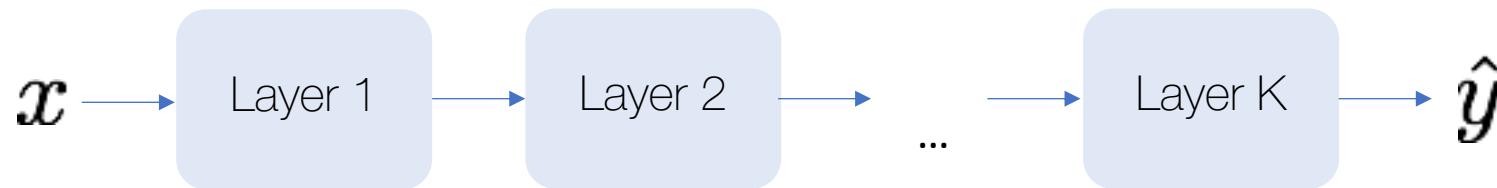


A Brief Intro to Deep Learning



We want to find $\{w_1, w_2, \dots, w_K\}$ s.t. $\hat{y}_i \approx y_i$

A Brief Intro to Deep Learning



We want to find $\{w_1, w_2, \dots, w_K\}$ s.t. $\hat{y}_i \approx y_i$

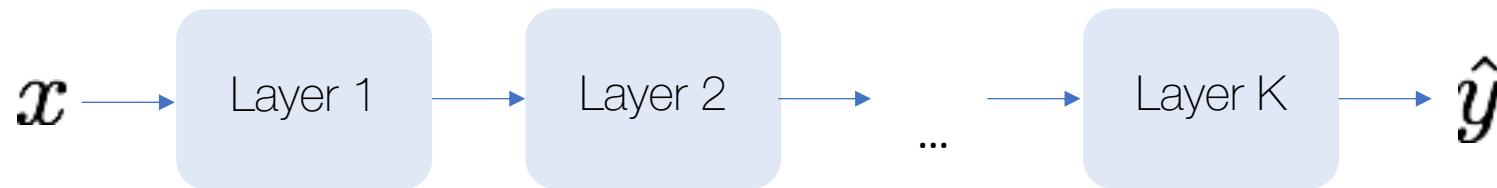
Loss functions

Quantify the error between prediction and ground truth

e.g.

$$l = \frac{1}{N} \sum_{i=0}^{N-1} (\hat{y}_i - y_i)^2$$

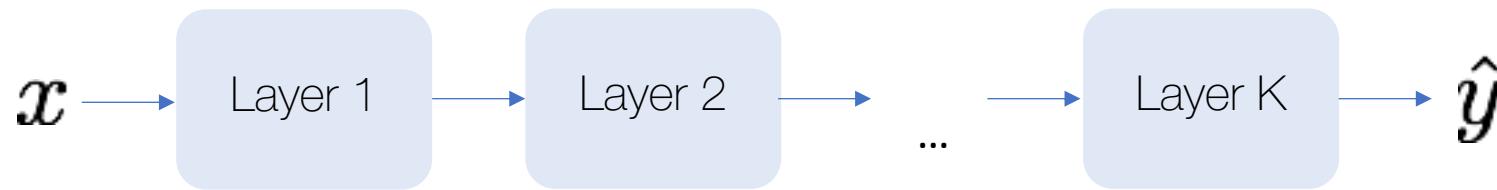
A Brief Intro to Deep Learning



Optimization

$$\min_{\{w_1, w_2, \dots, w_K\}} l = \min_{\{w_1, w_2, \dots, w_K\}} \frac{1}{N} \sum_{i=0}^{N-1} (\hat{y}_i - y_i)^2$$

A Brief Intro to Deep Learning



Optimization: Gradient Descent

At t=0, $w_j^{(0)} = 0$

At t-th iteration:

$$w_j^{(t)} = w_j^{(t-1)} - \eta \frac{dl}{dw_j}$$

Learning rate

Gradient

A Brief Intro to Deep Learning



Remember that
 $x = in_1$
 $\hat{y} = out_K$
 $out_j = in_{j+1}$

Gradient Computation: Chain Rule

$$\frac{dl}{dw_j} = \frac{dl}{d\hat{y}} \cdot \frac{dout_K}{din_K} \cdot \frac{dout_{K-1}}{din_{K-1}} \cdot \dots \cdot \frac{dout_j}{dw_j}$$

Visual Recognition Tasks



Digit recognition



Image Classific.



Semantic Segmentation



Object Detection



Instance Segmentation

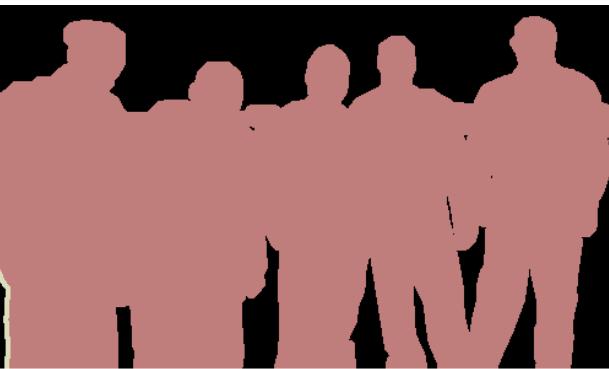


Increasing complexity

Content

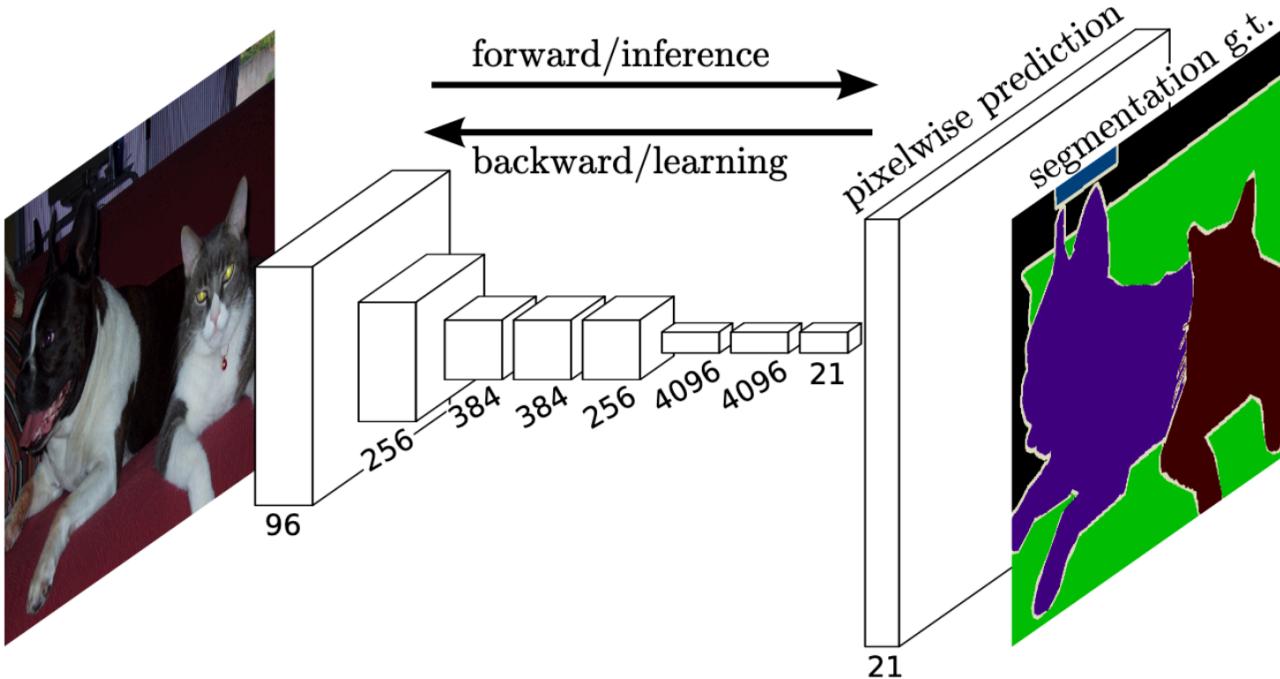
- A Brief Intro
- Visual Recognition: Problem Definitions & Tasks
 - Semantic Segmentation
 - FCN
 - Object Detection
 - Fast(er) R-CNN
 - Instance Segmentation & Pose Prediction
 - Mask R-CNN
- ImageNet in 1H: Important Lessons
- Useful Tips

Semantic Segmentation



- Given an input image I , mark each pixel in I with a class label
- Unlike image classification, predictions are cast for every pixel
- Metrics:
 - Mean Accuracy
 - Mean Intersection over Union

FCN

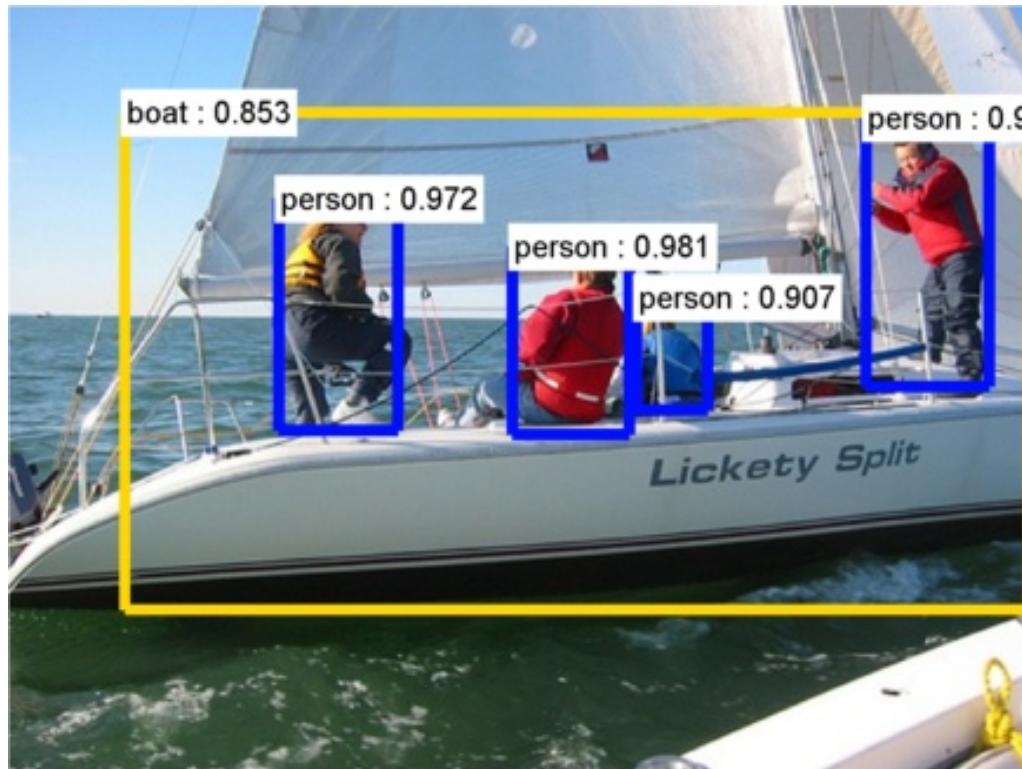


- Input image I of size $H \times W \times 3$
- Output logits of size $H \times W \times C$, where C is the number of object categories in the dataset.

FCN: Losses

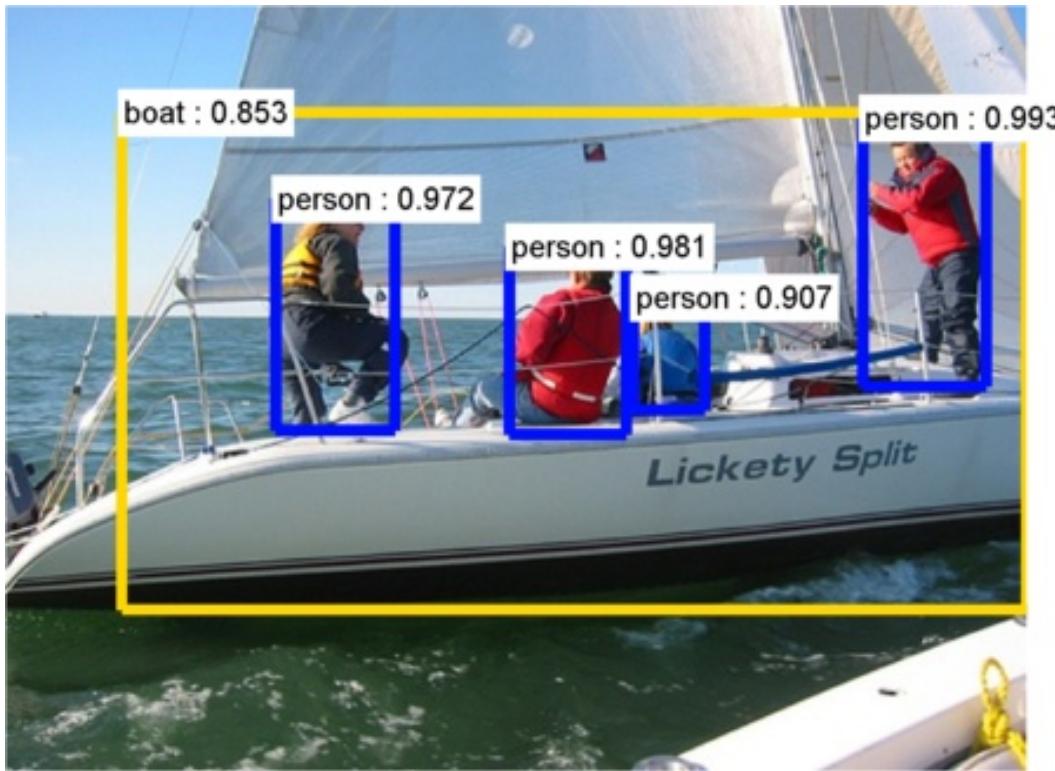
- Assume $J \in \mathbb{R}^{H \times W \times C}$ are the logits predicted by FCN & $J_{gt} \in [0, C - 1]^{H \times W}$ is the ground truth
- Binary classification vs. Softmax

Object Detection



What
Where

Object Detection

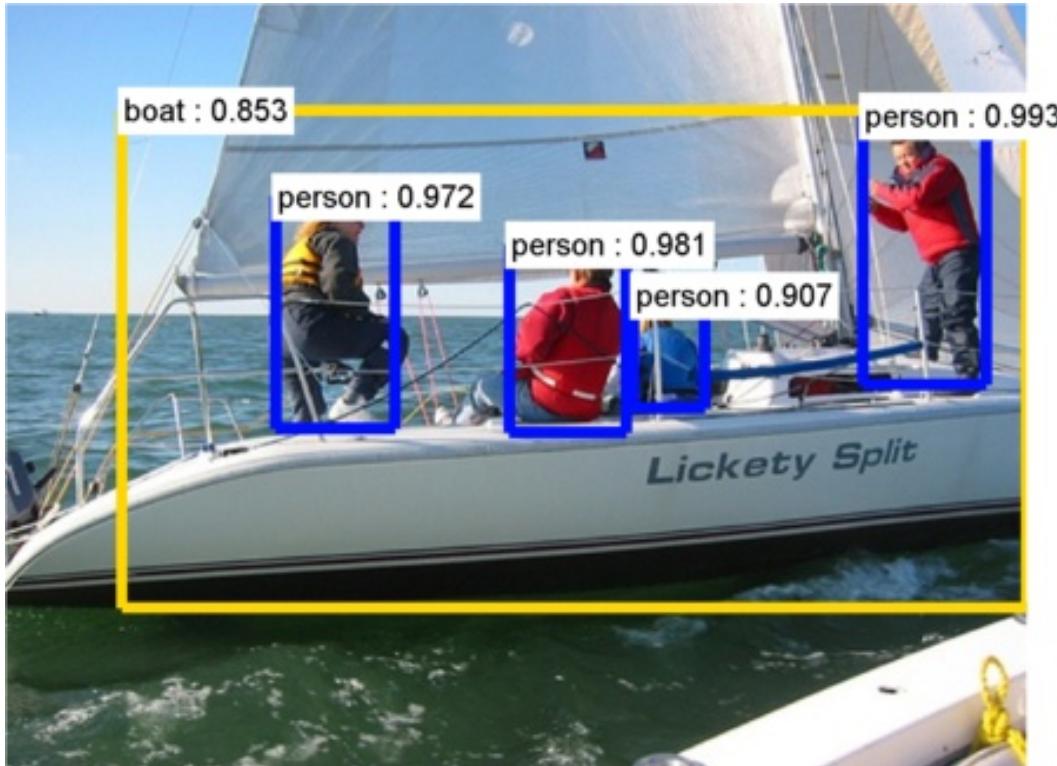


Object Detection



Image Classification

Object Detection



Task

- Assume C object classes
- Localize and classify all objects in an image

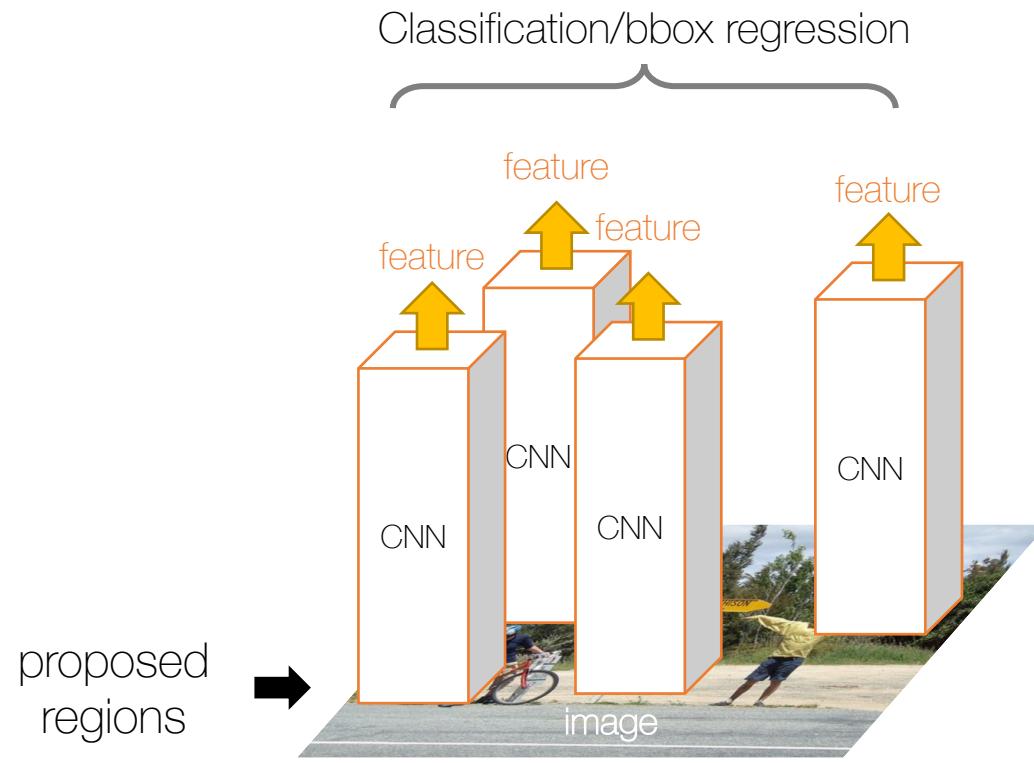
True Positives

- A detection b is a true positive if $ov(b, gt) > 0.5$ and $class(b) = class(gt)$

False Positives

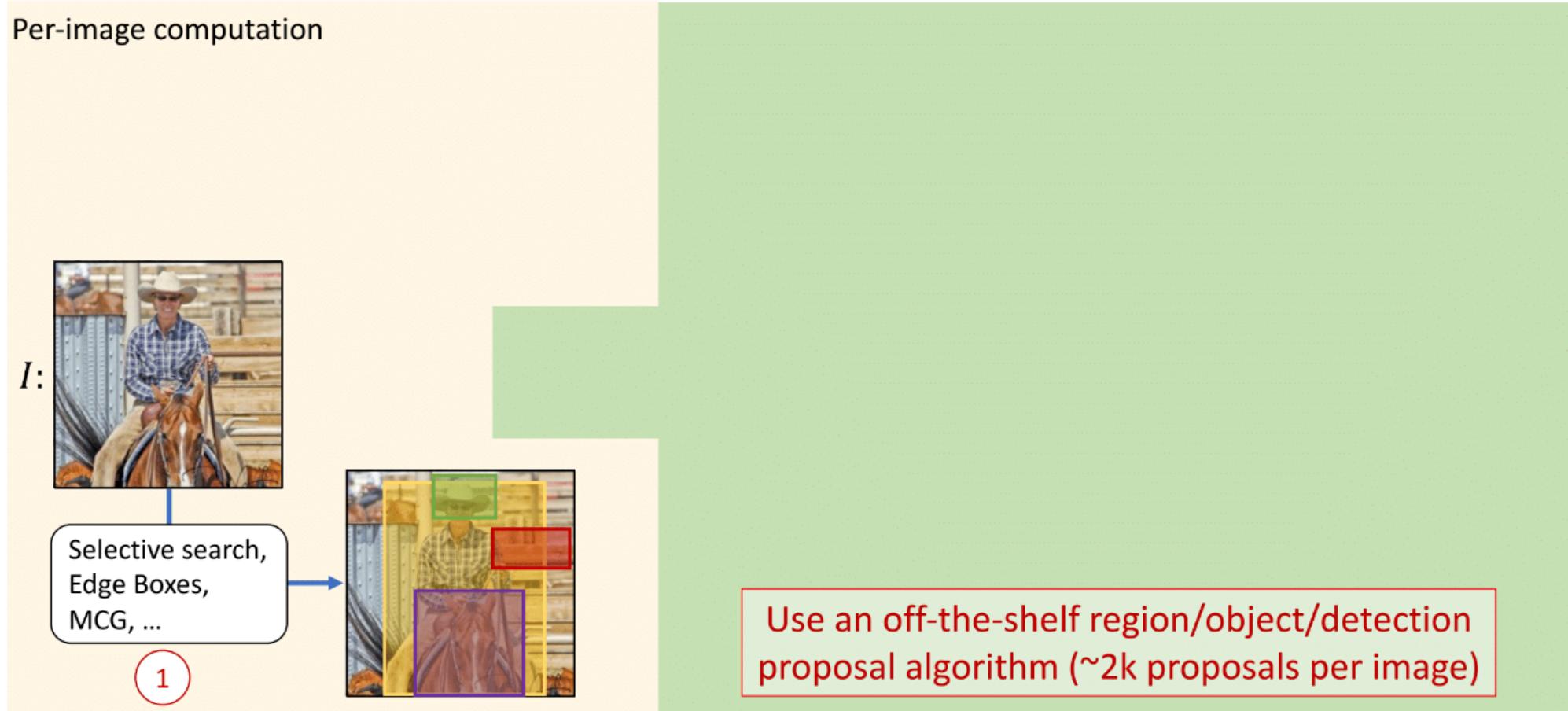
- Mislocalized and misclassified detections
- Duplicate detections are penalized

Object Detection Frameworks: R-CNN

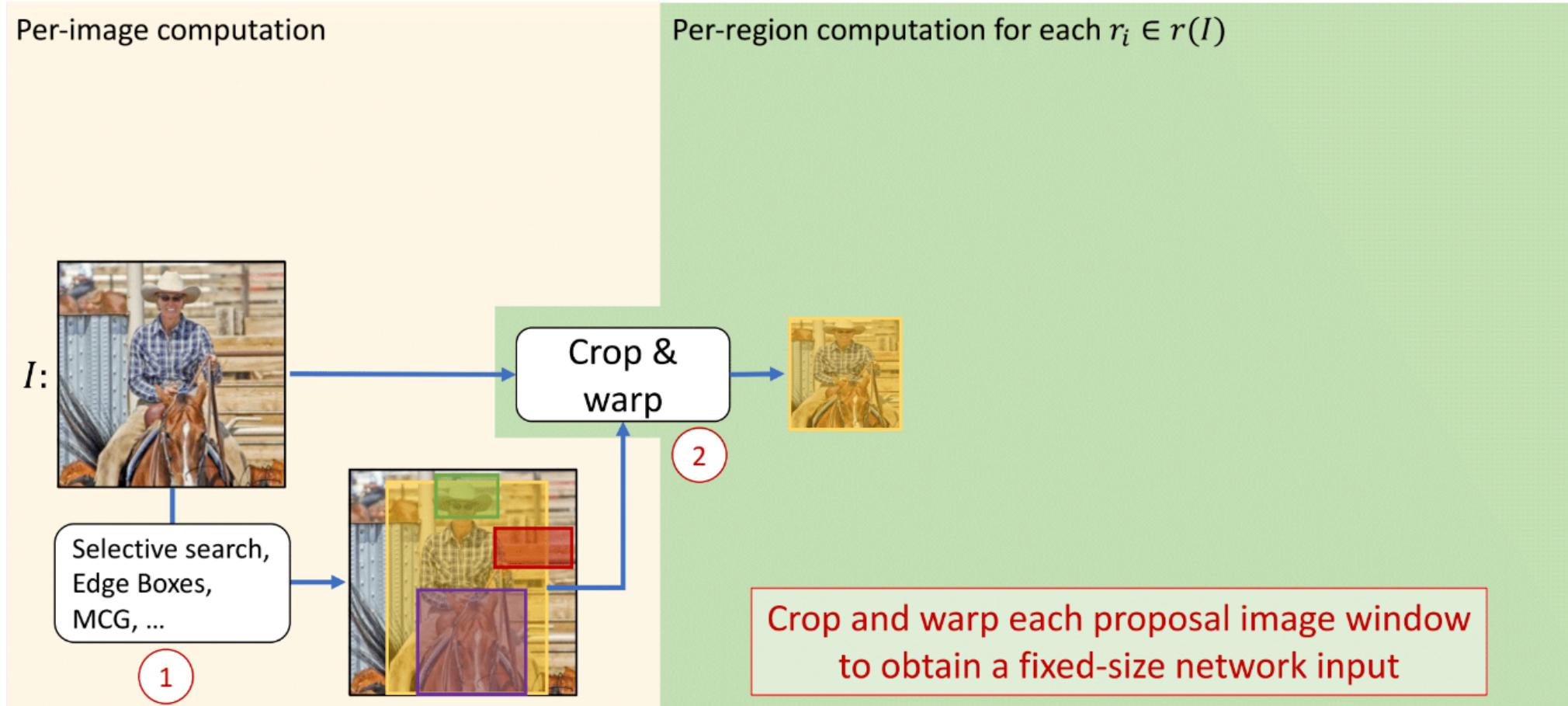


R-CNN

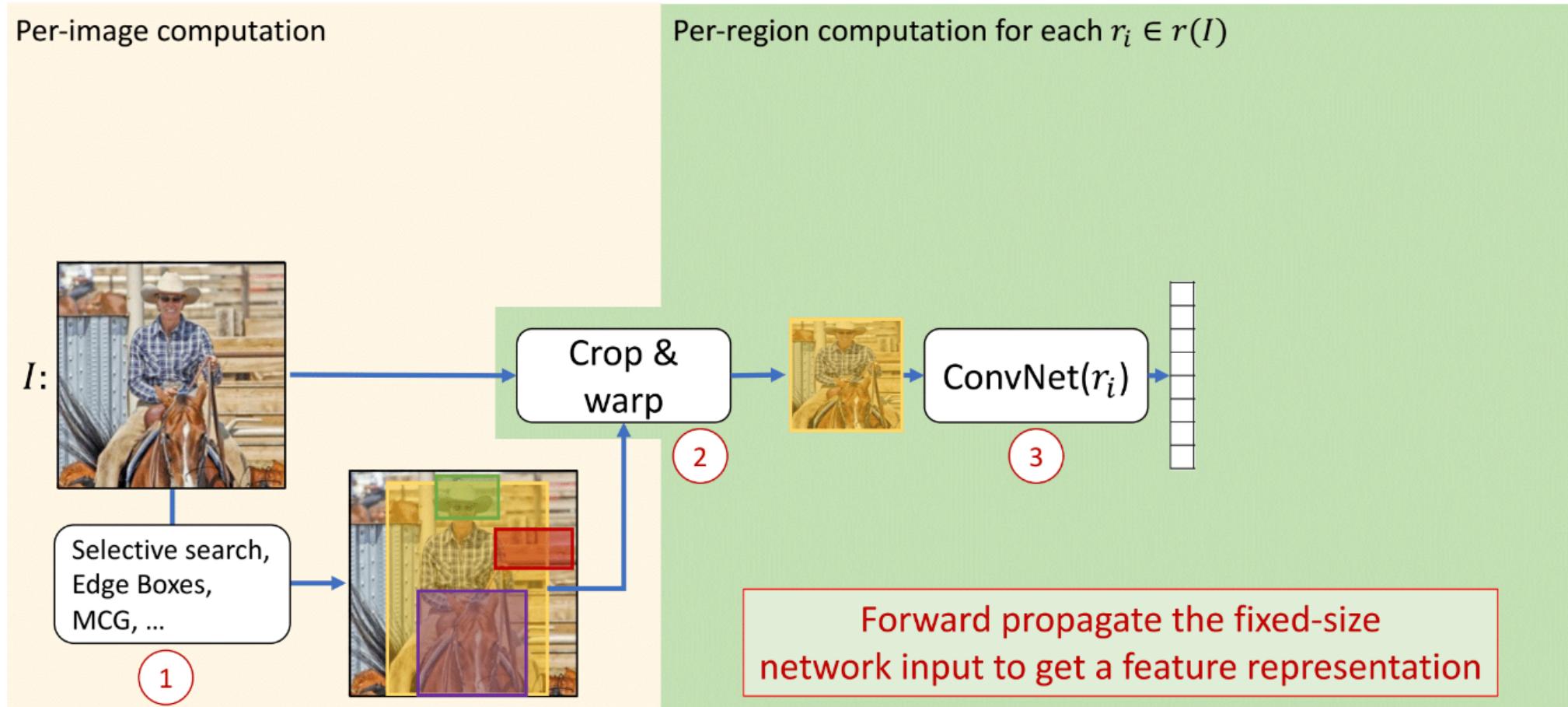
Object Detection Frameworks: R-CNN



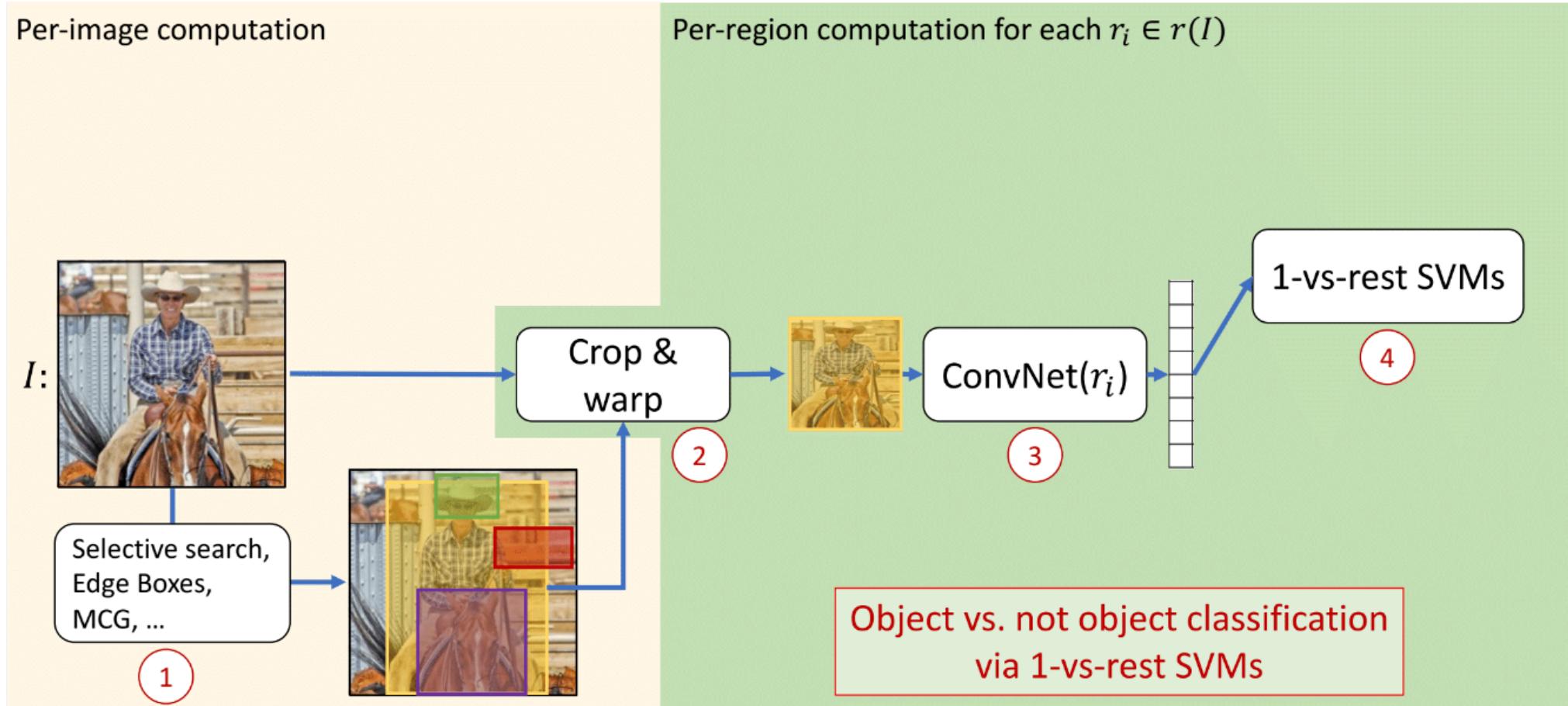
Object Detection Frameworks: R-CNN



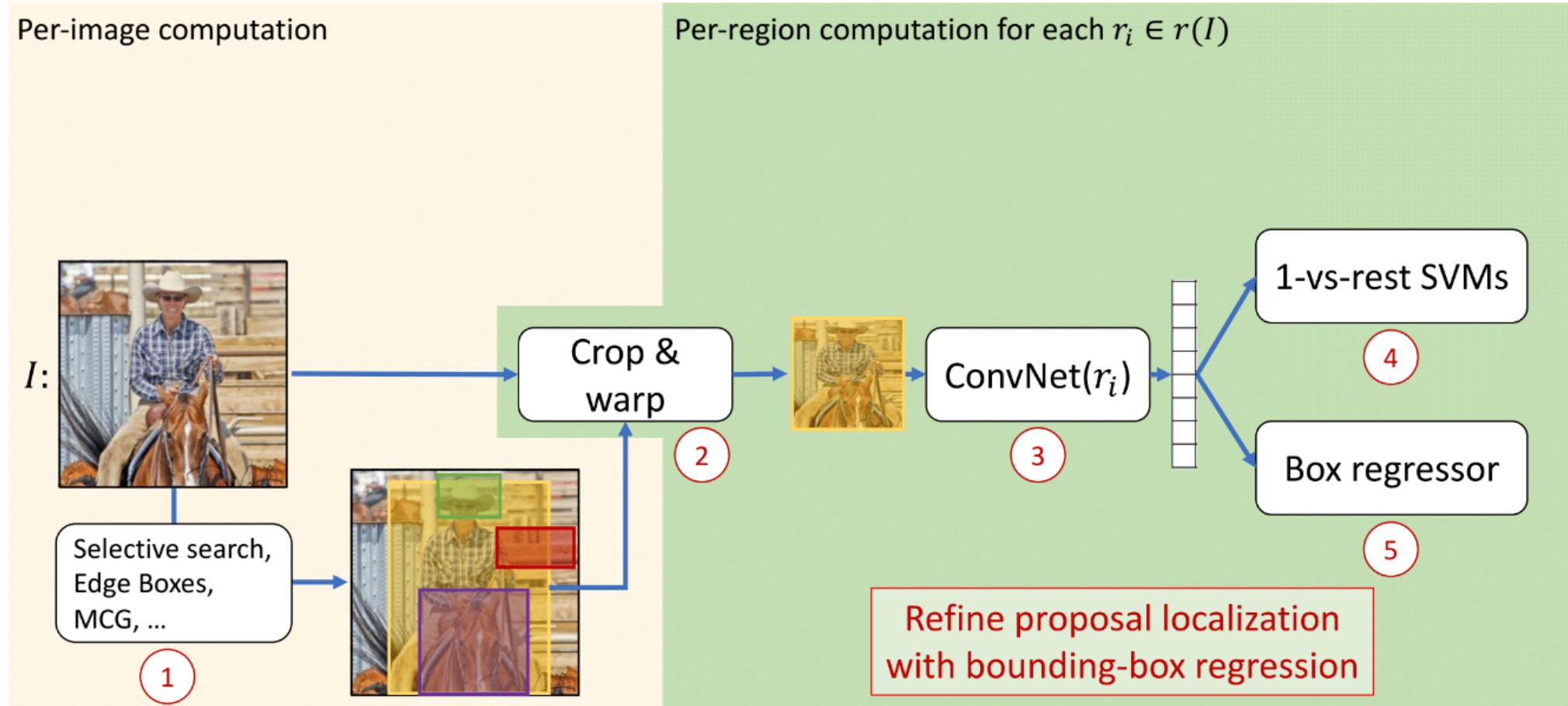
Object Detection Frameworks: R-CNN



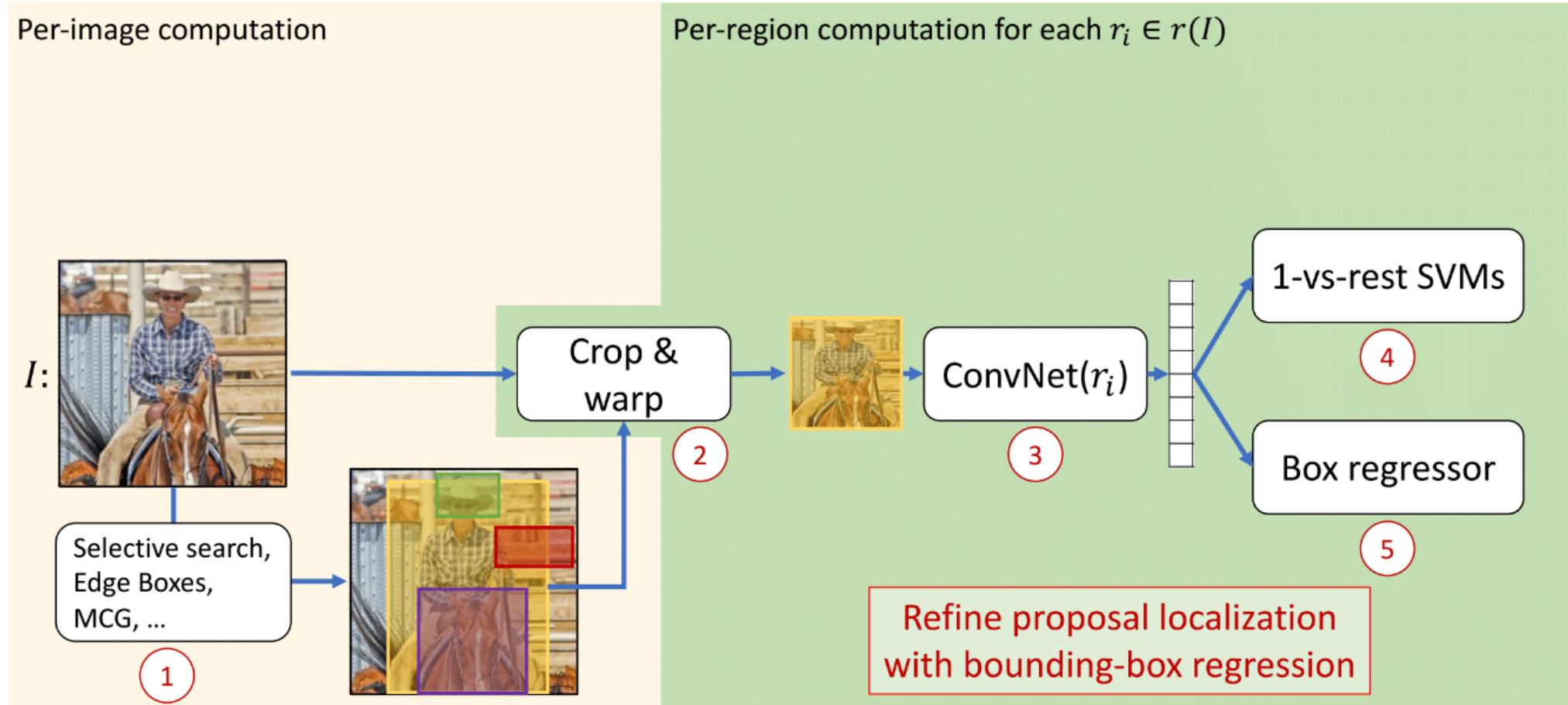
Object Detection Frameworks: R-CNN



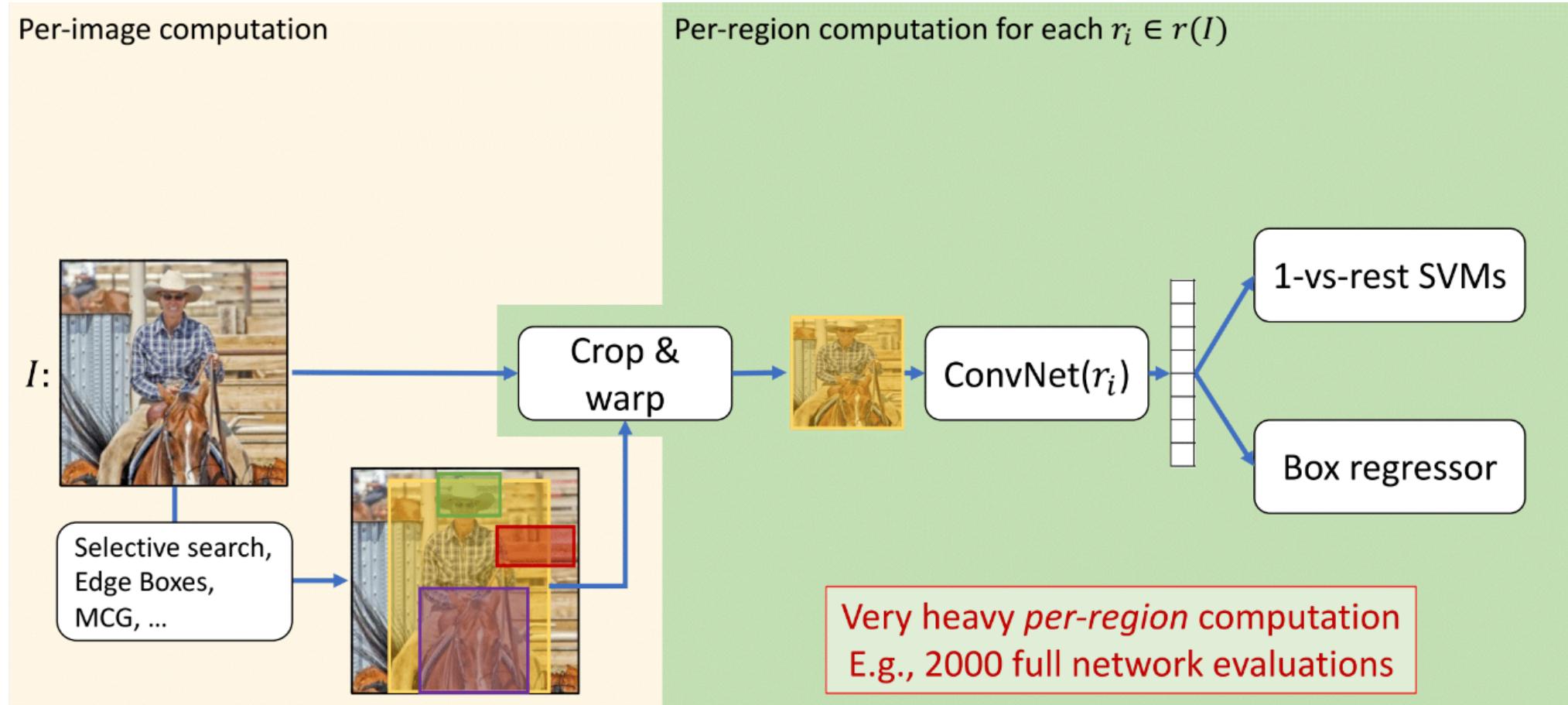
Object Detection Frameworks: R-CNN



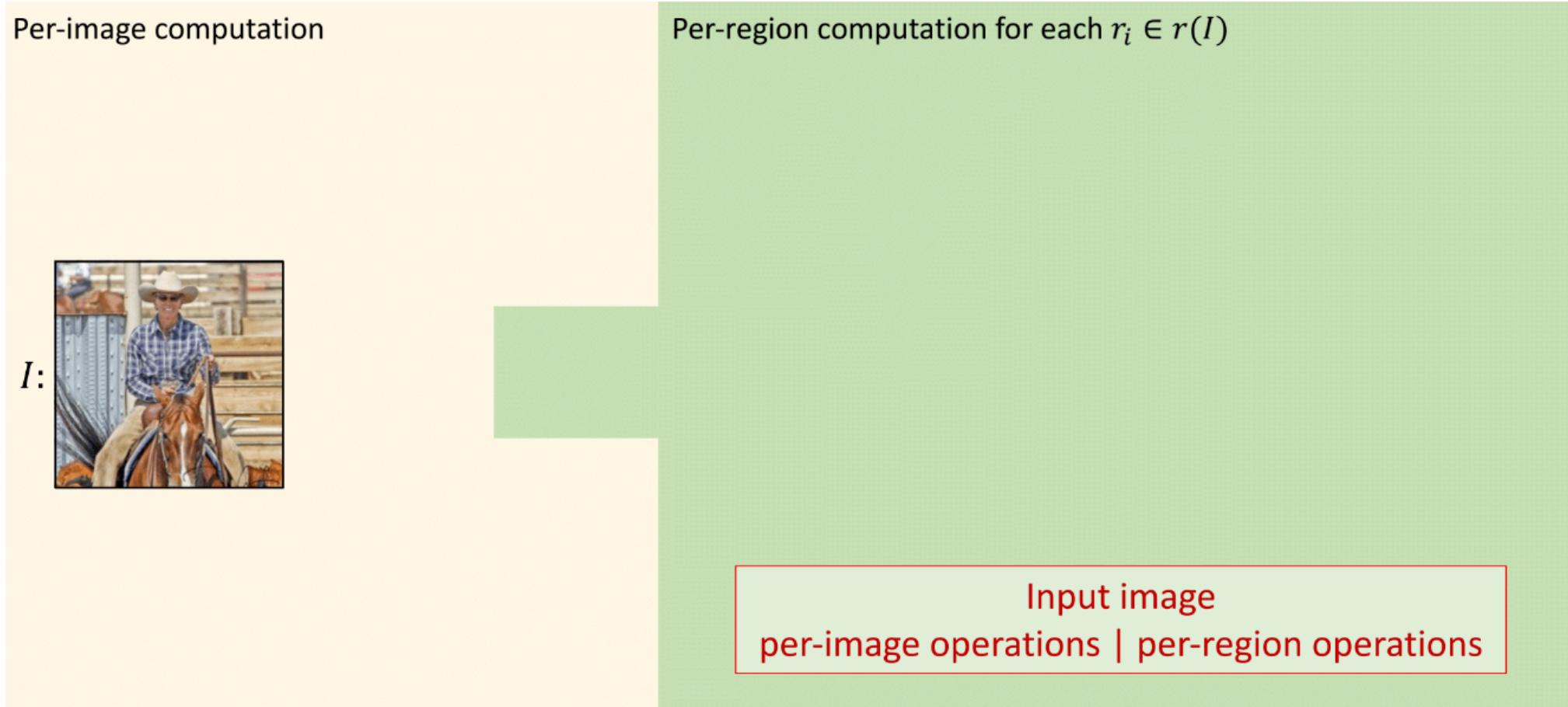
Object Detection Frameworks: R-CNN



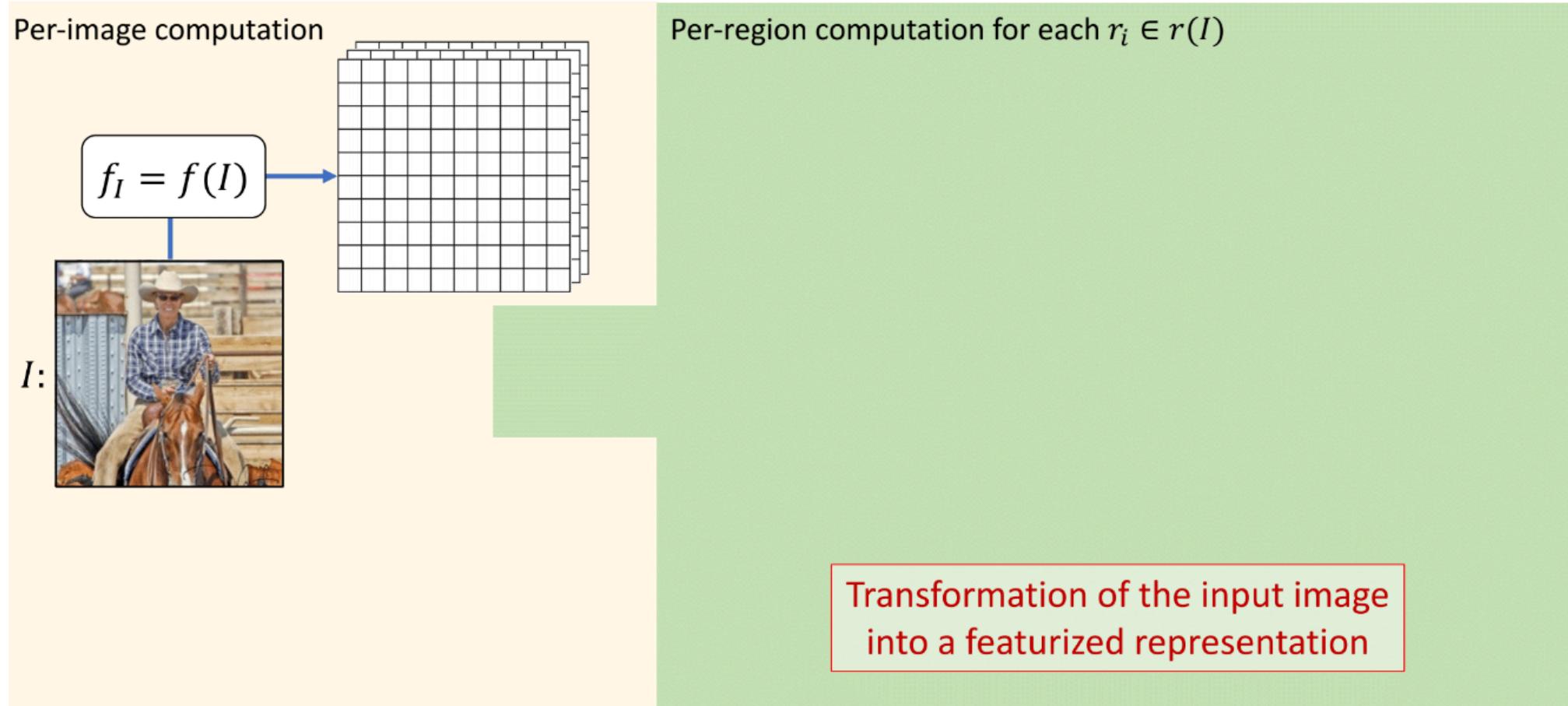
Object Detection Frameworks: R-CNN



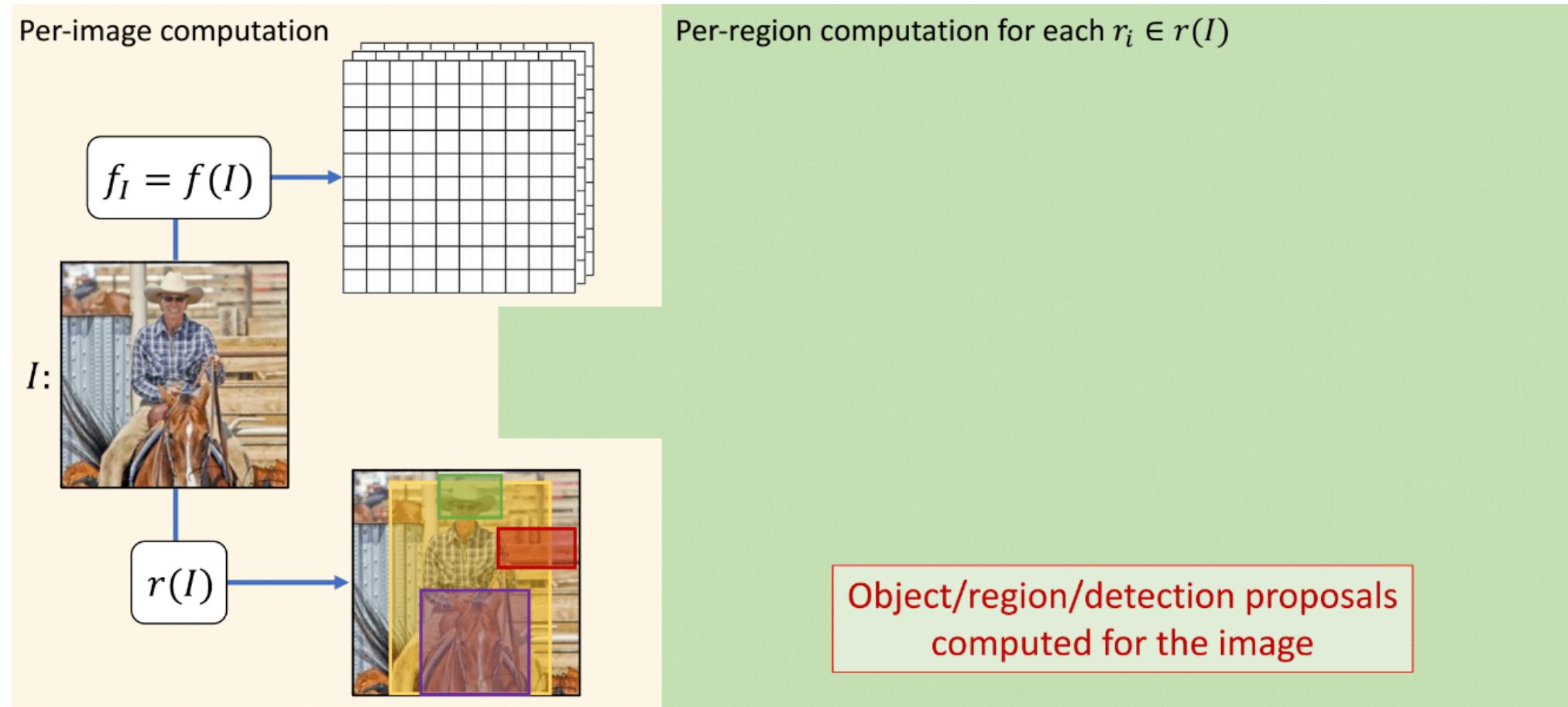
Object Detection Frameworks



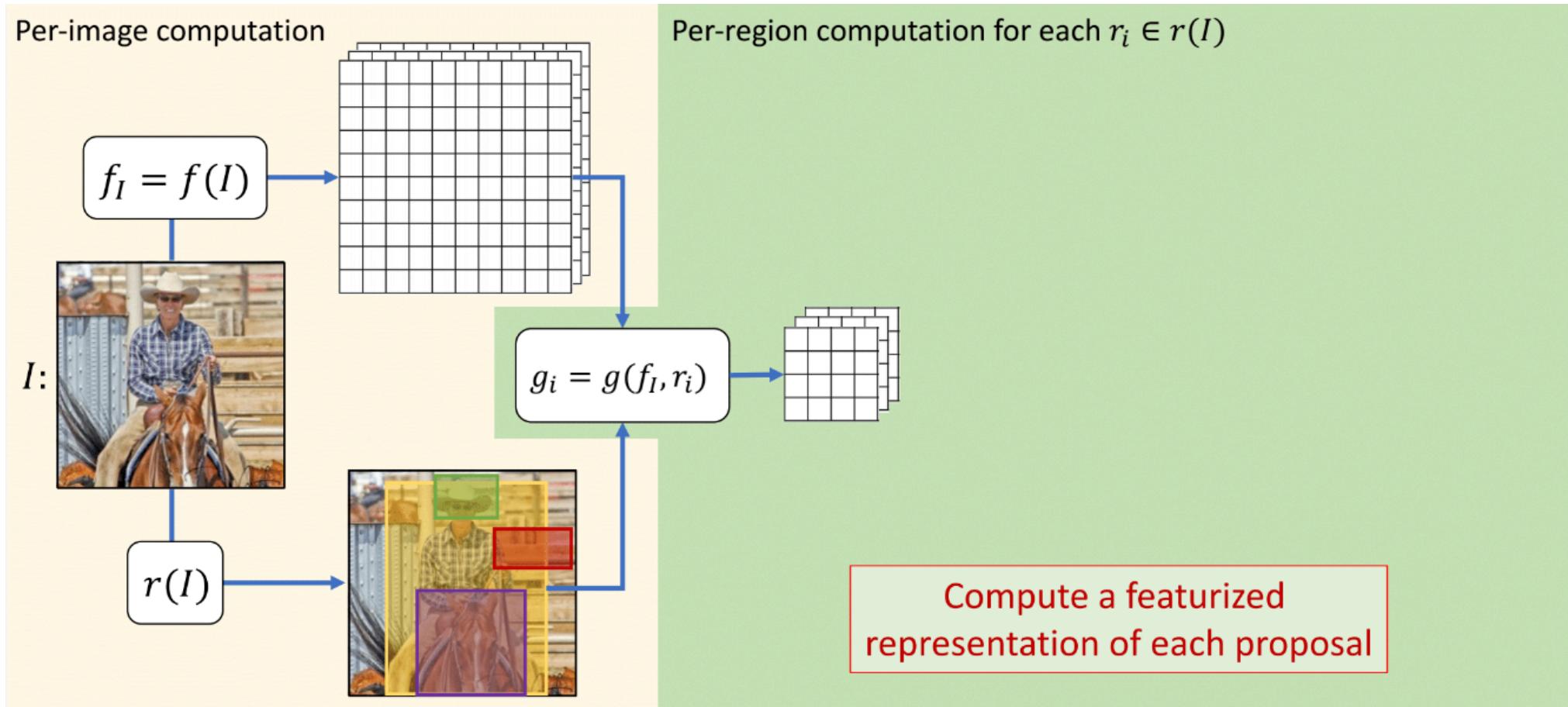
Object Detection Frameworks



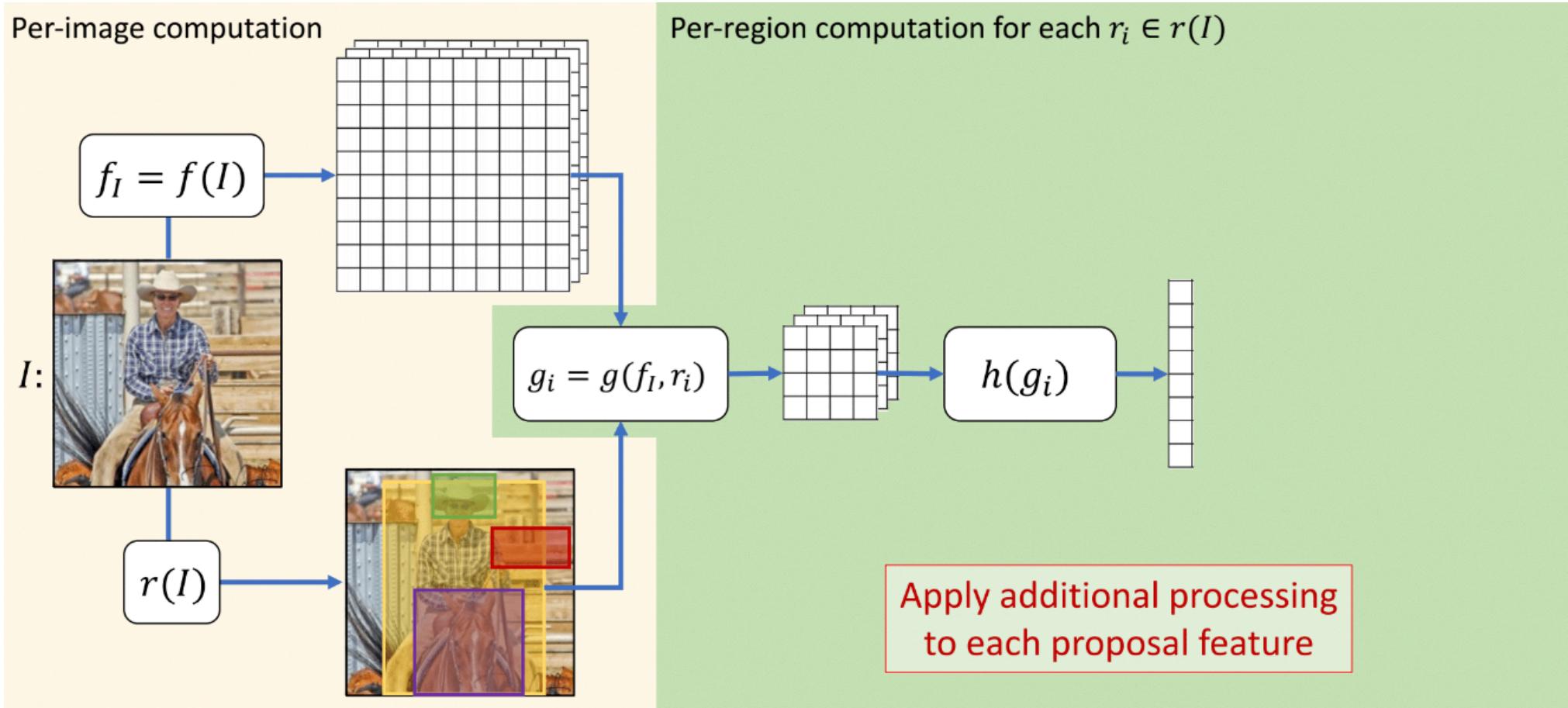
Object Detection Frameworks



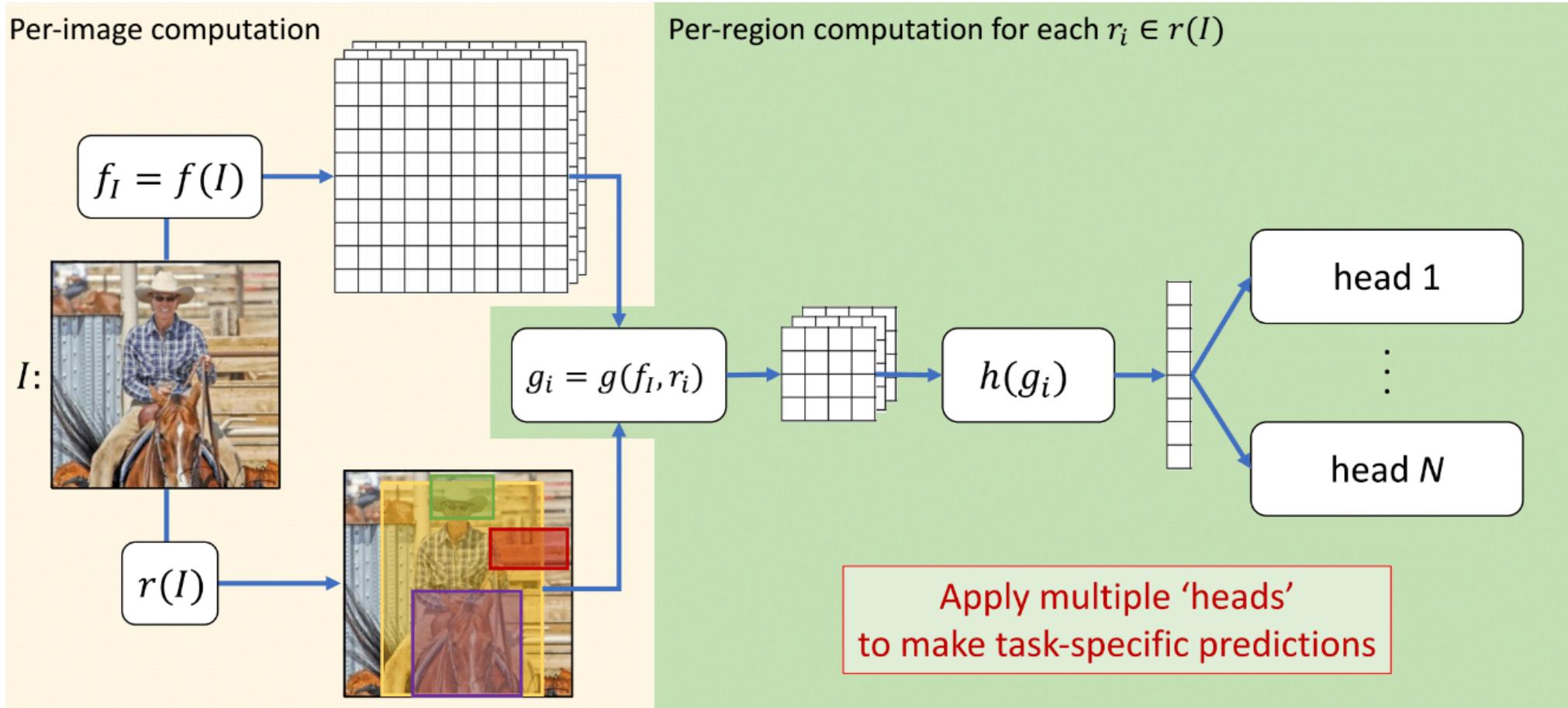
Object Detection Frameworks



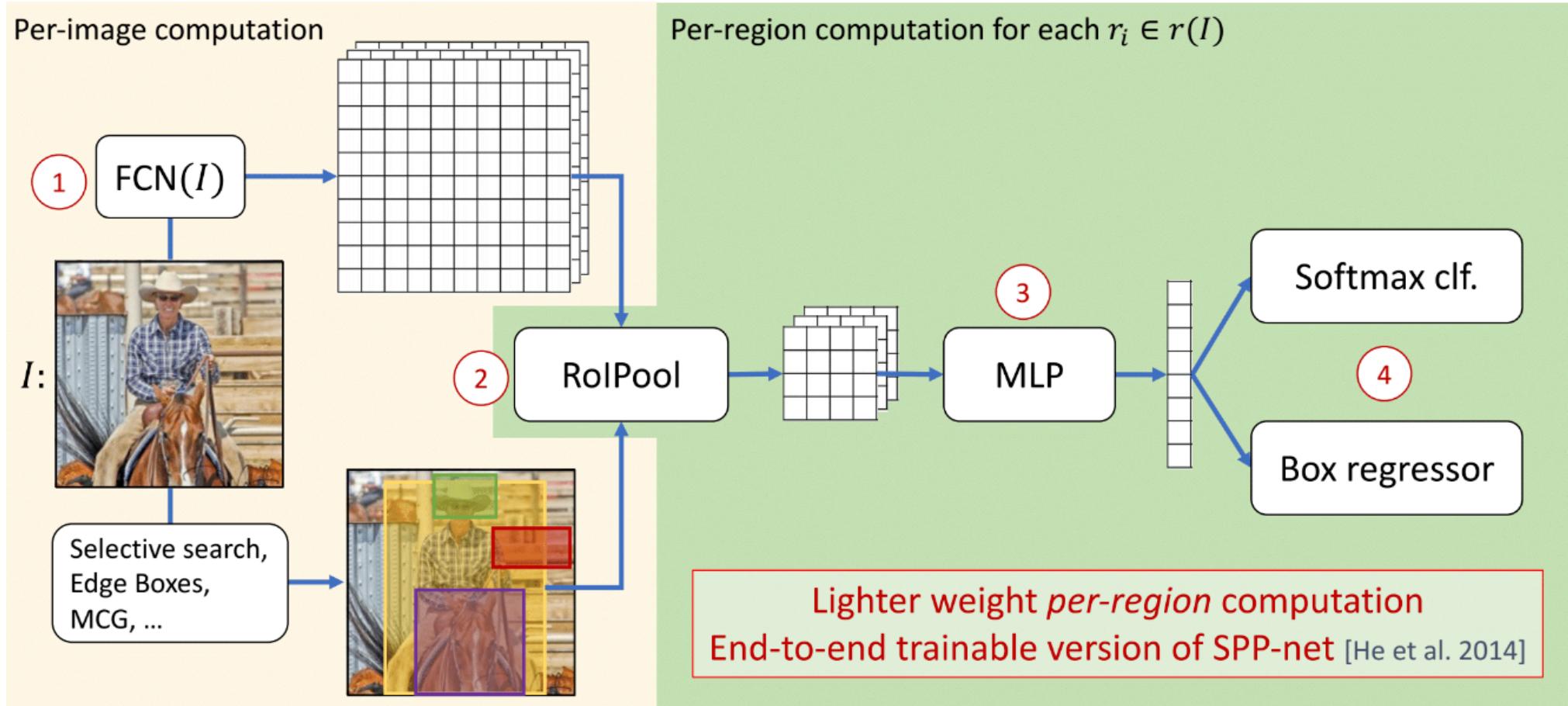
Object Detection Frameworks



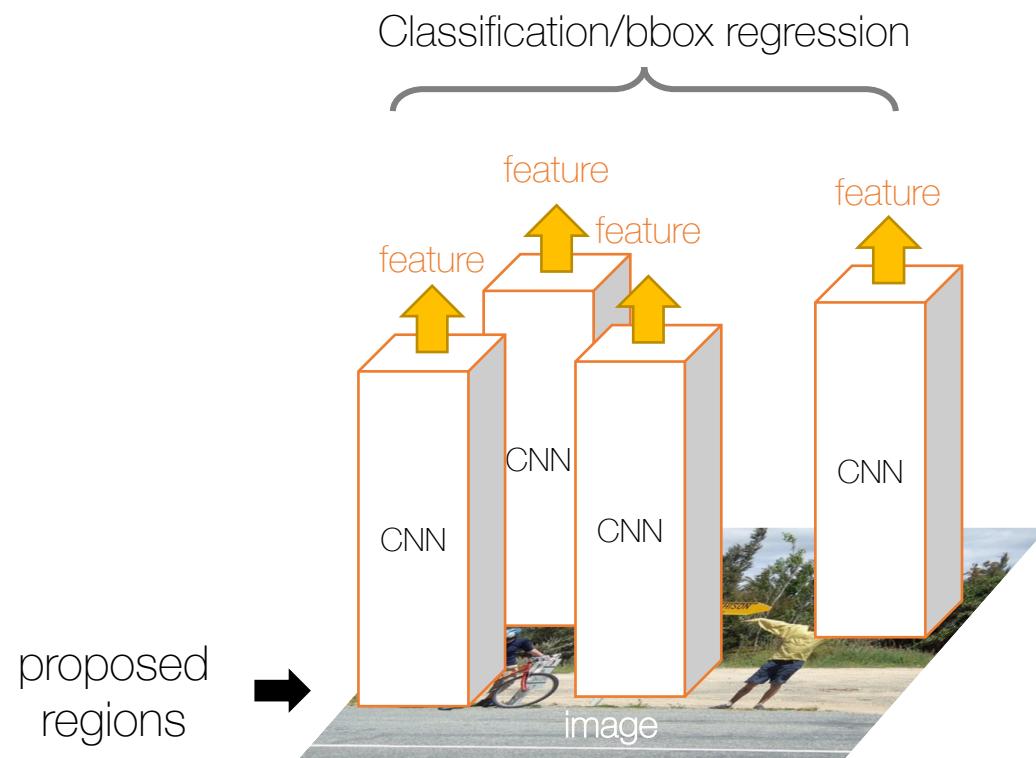
Object Detection Frameworks



Object Detection Frameworks: Fast R-CNN

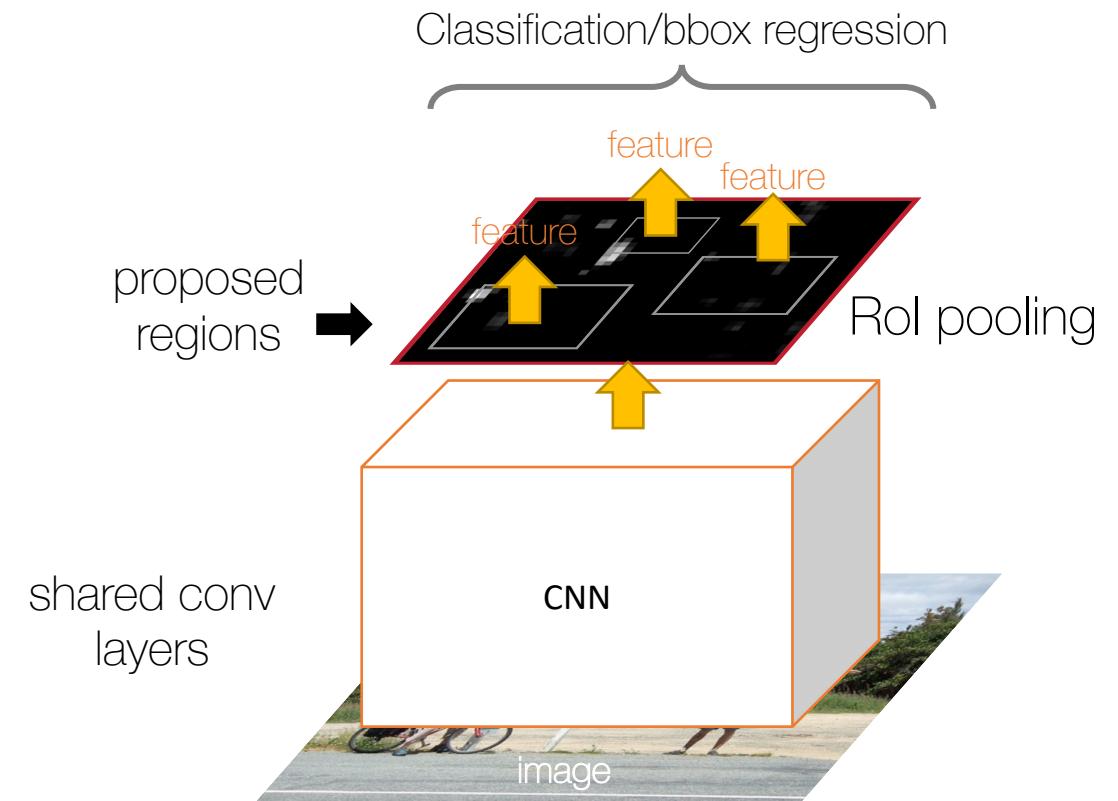


Object Detection Frameworks



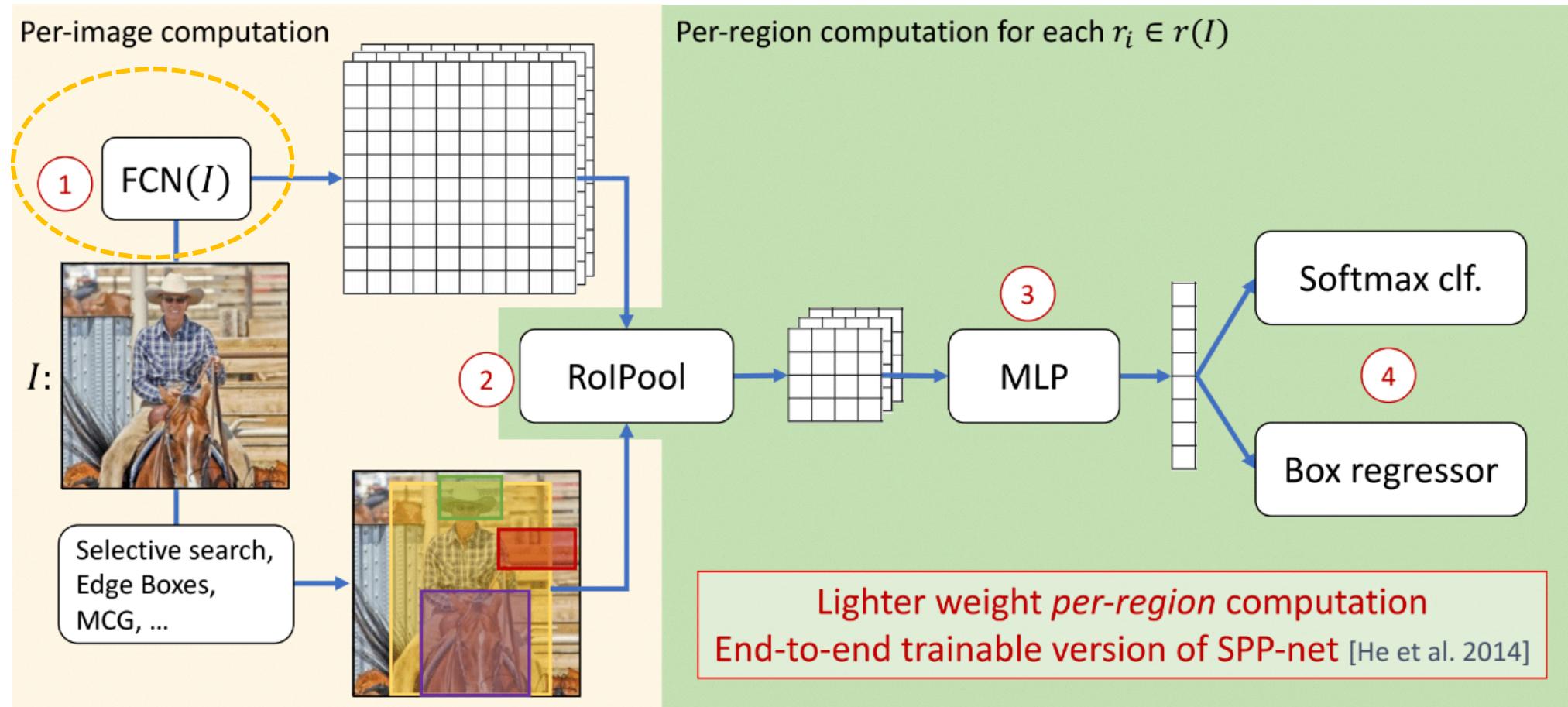
R-CNN

Girshick et al, CVPR 2014 & ICCV 2015

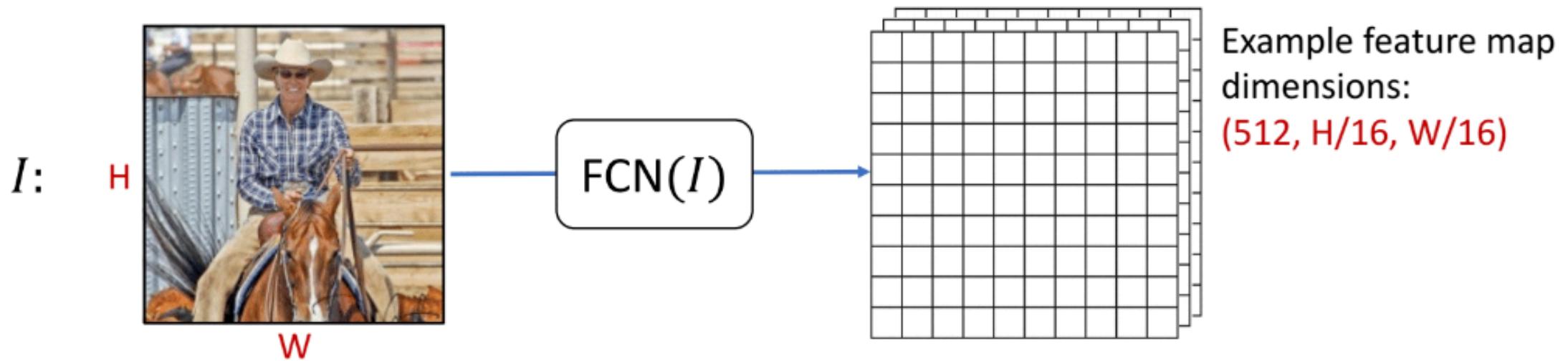


Fast R-CNN

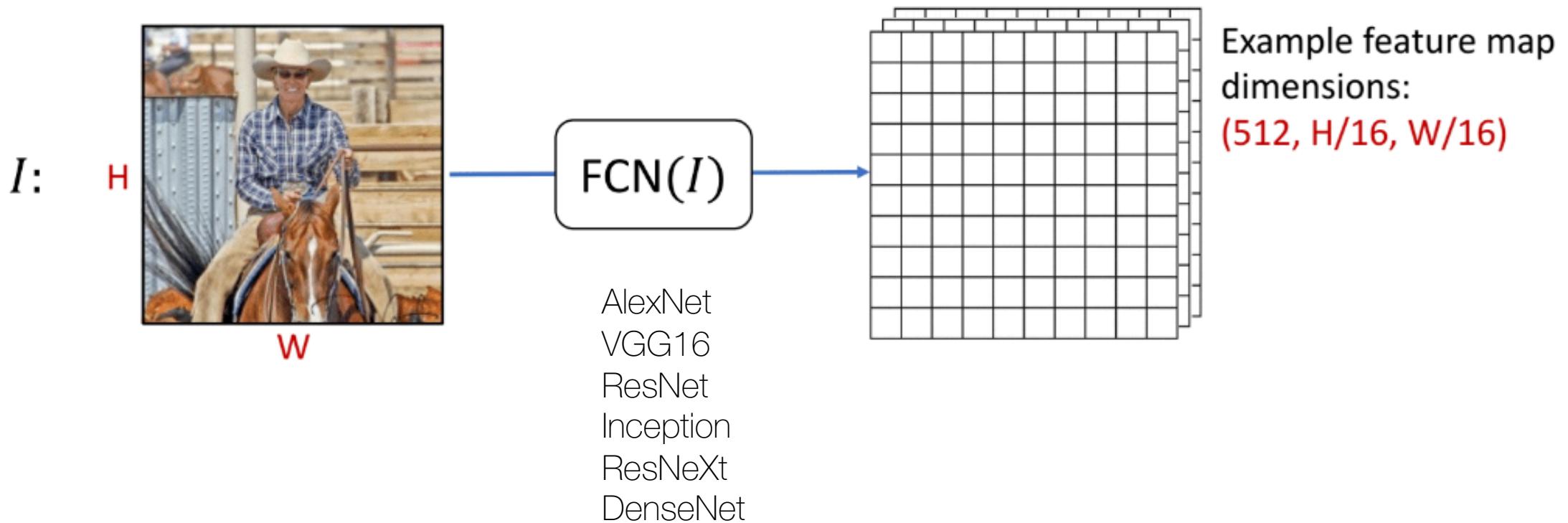
Backbone Networks



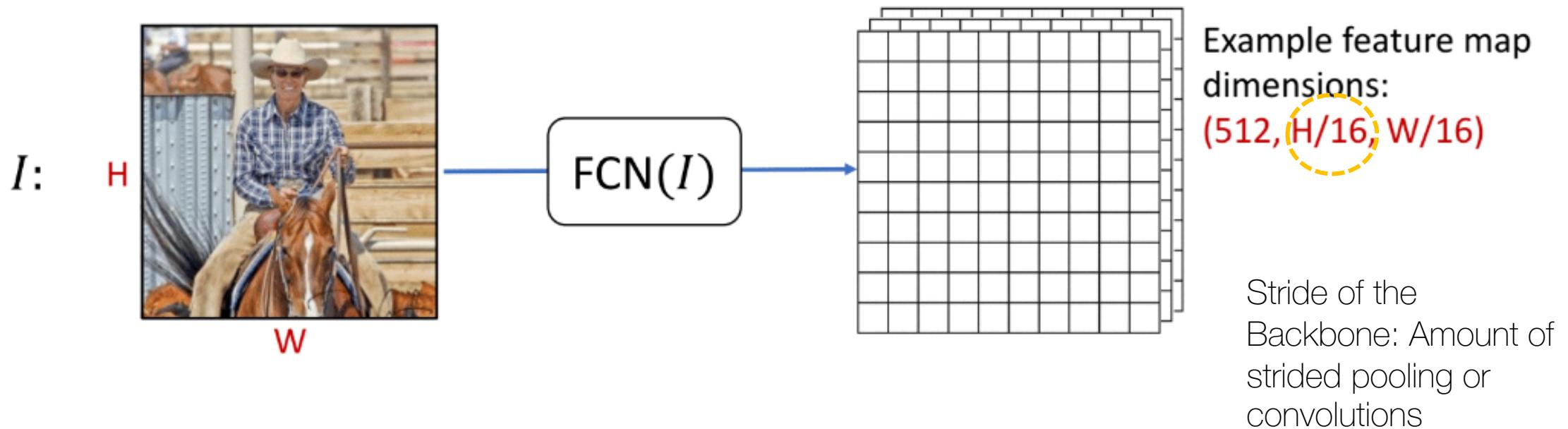
Backbone Networks



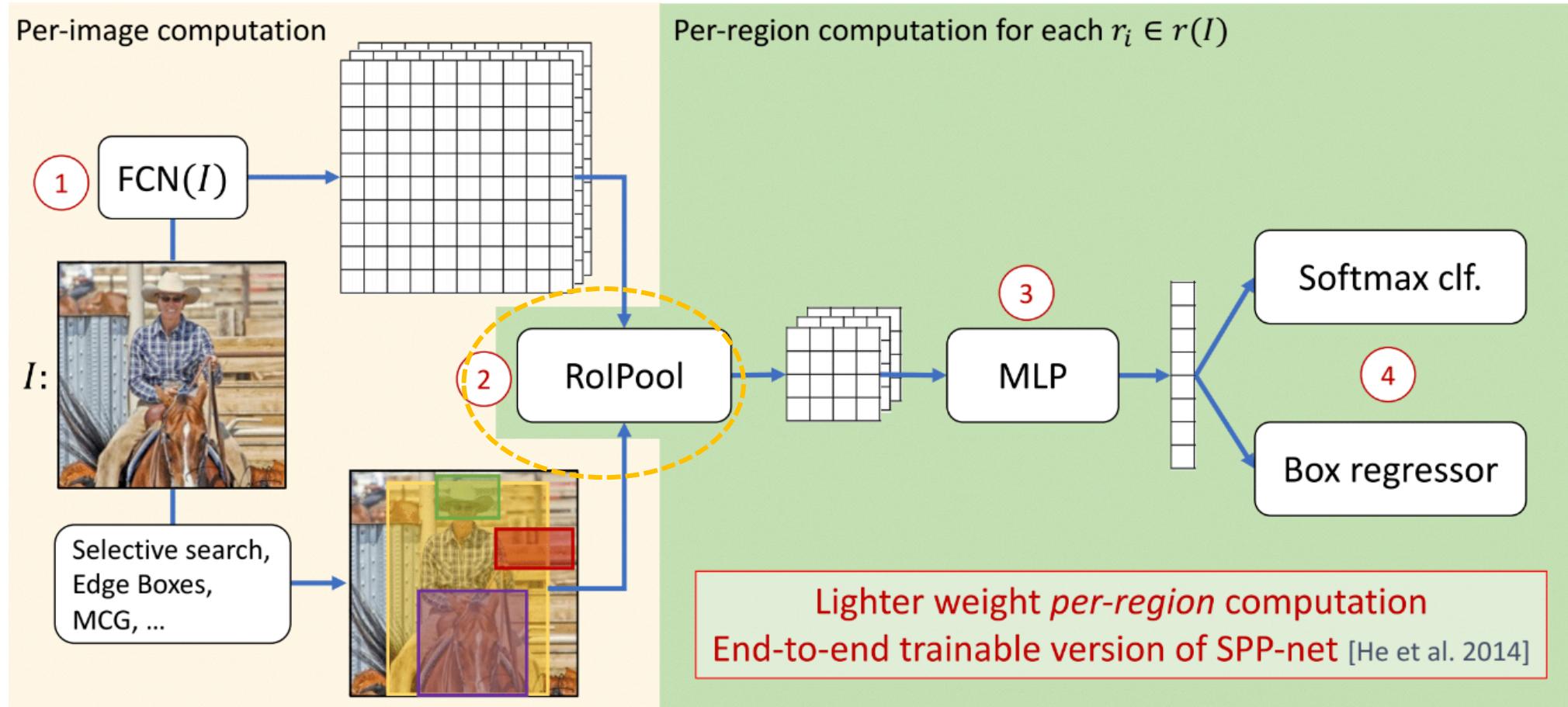
Backbone Networks



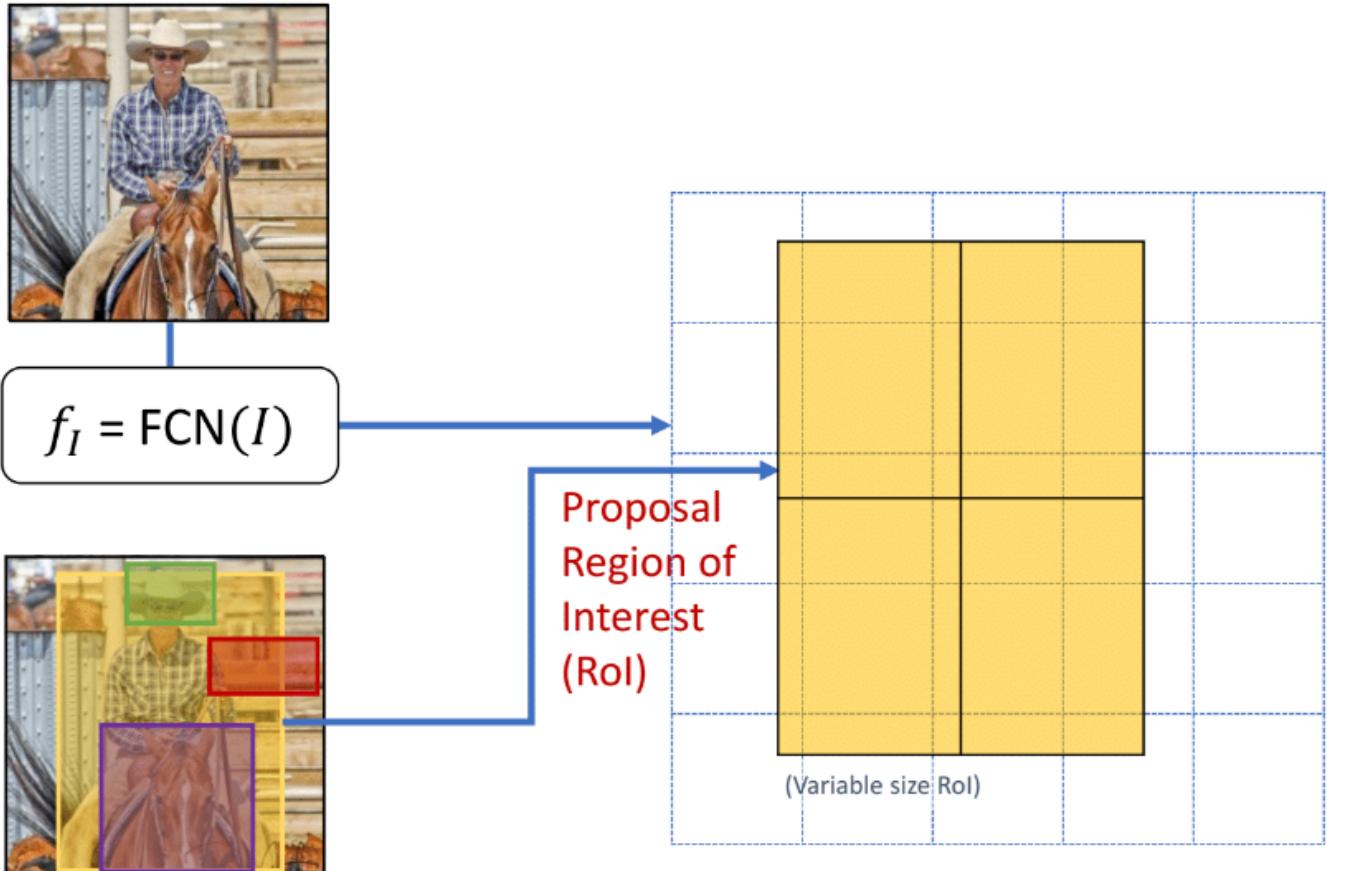
Backbone Networks



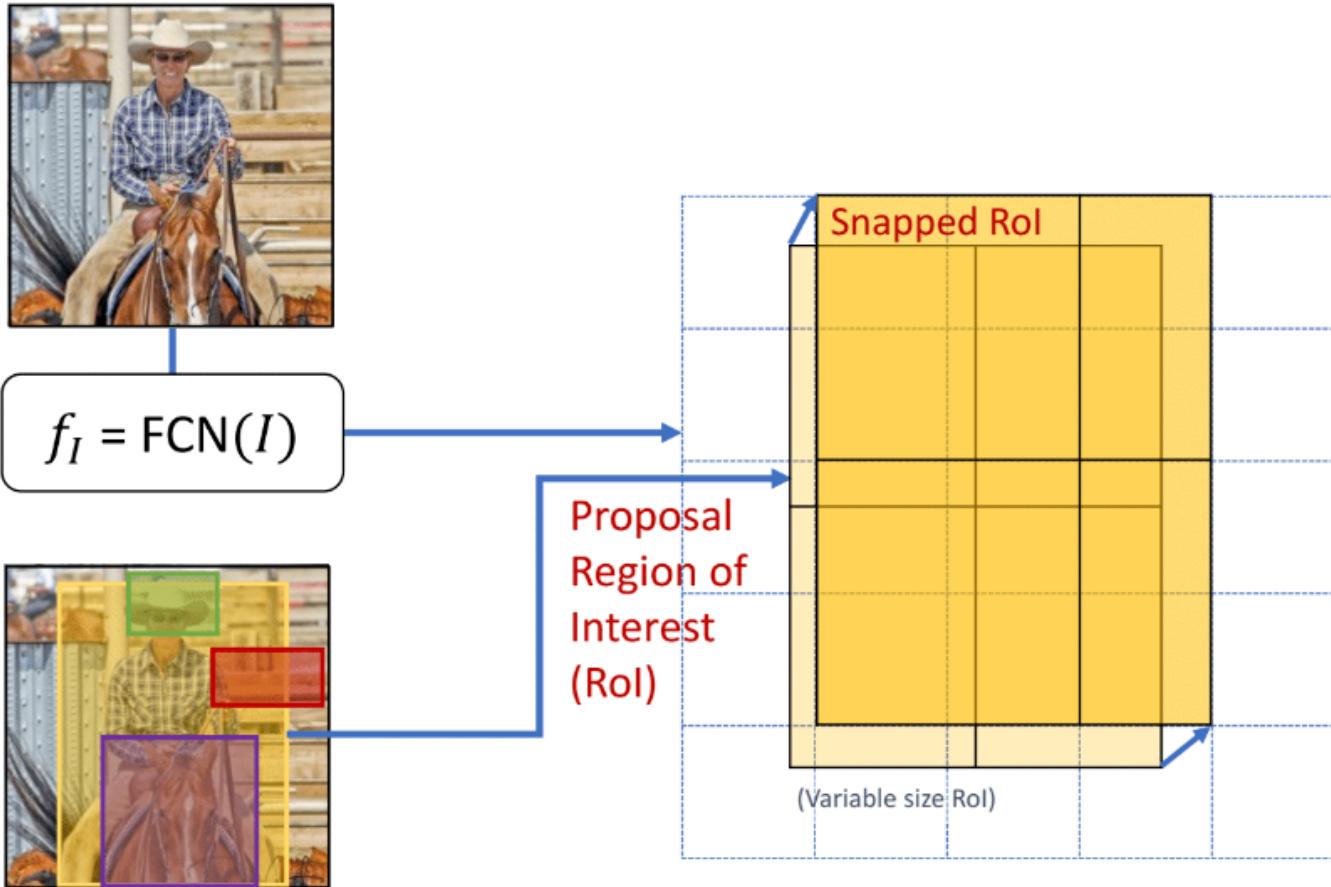
RoIPool



RoIPool on Each Proposal

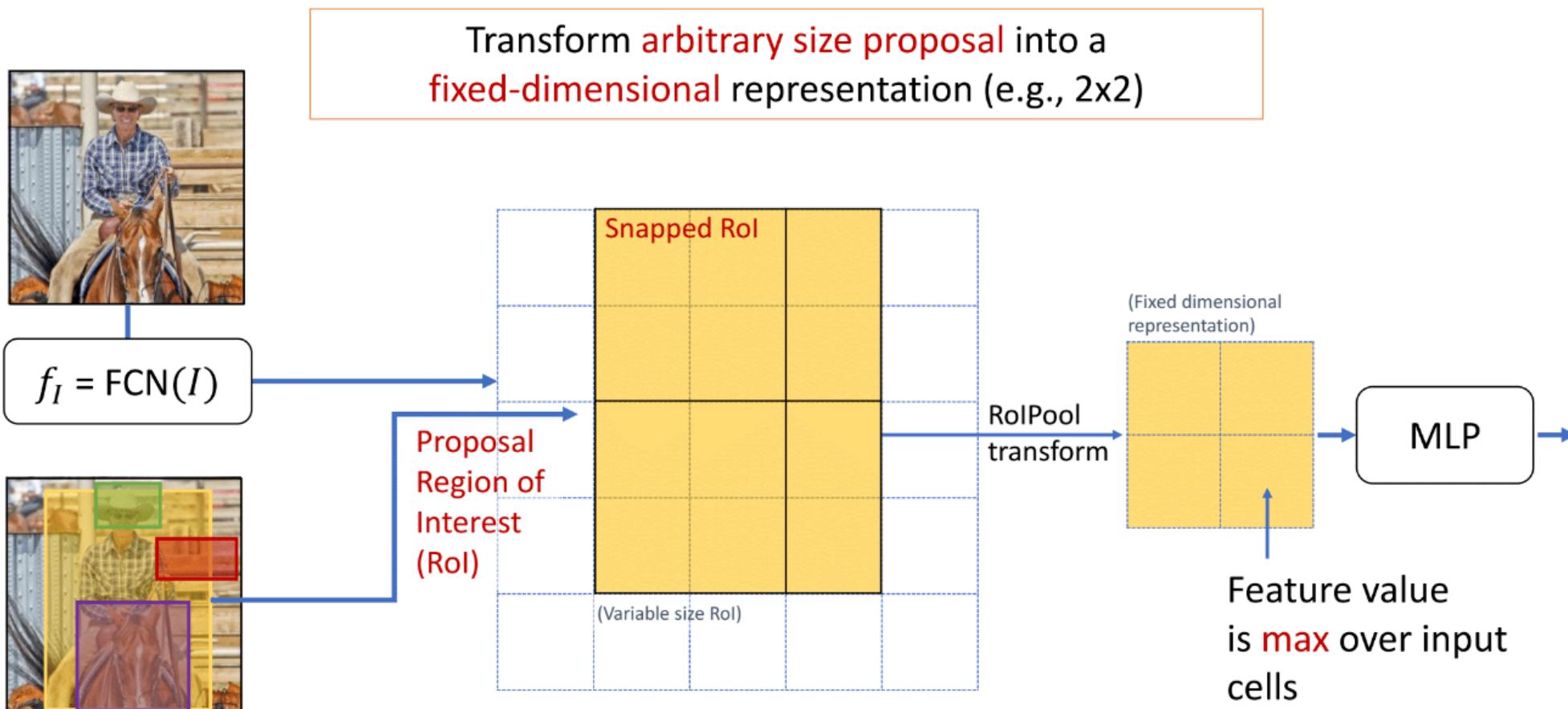


RoIPool on Each Proposal

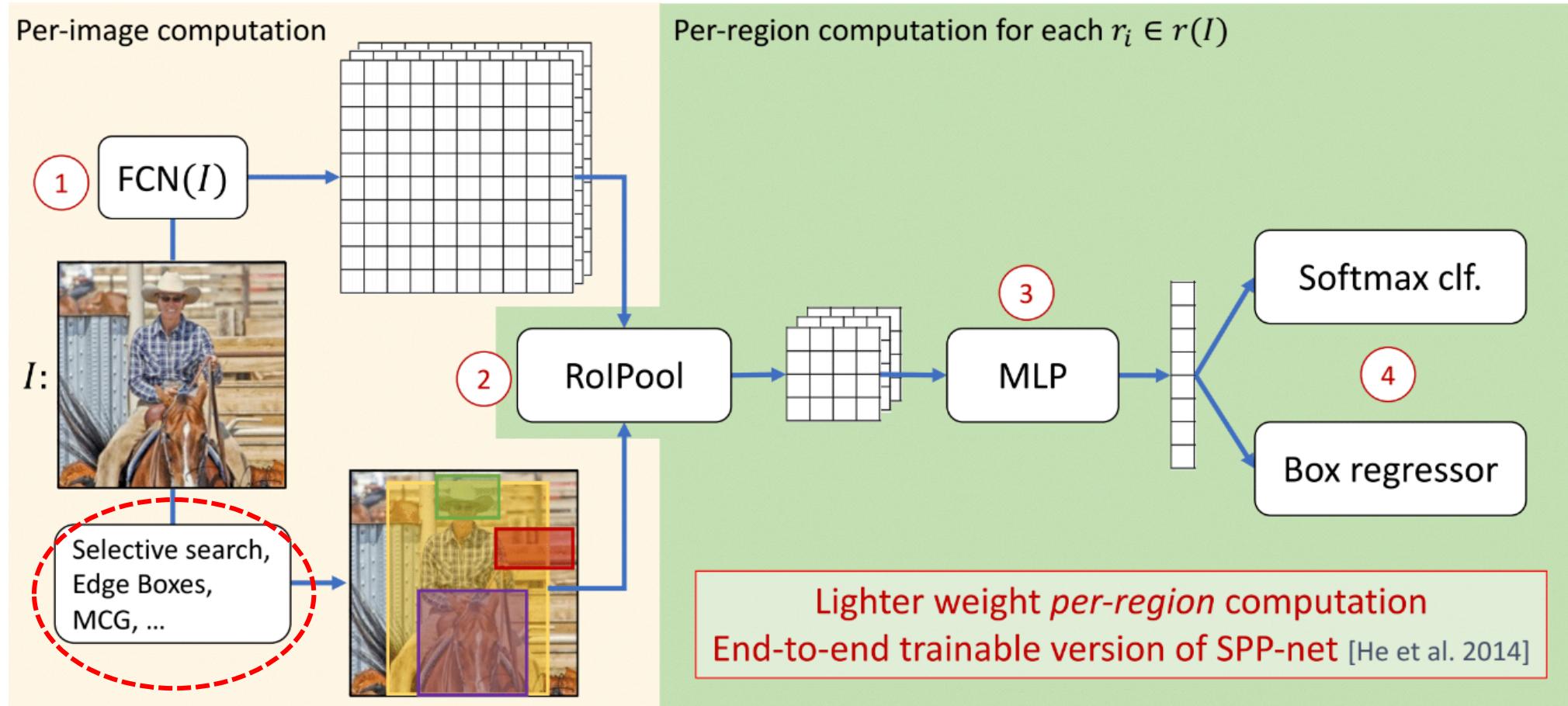


Innovation in SPP-Net [He et al. 2014]

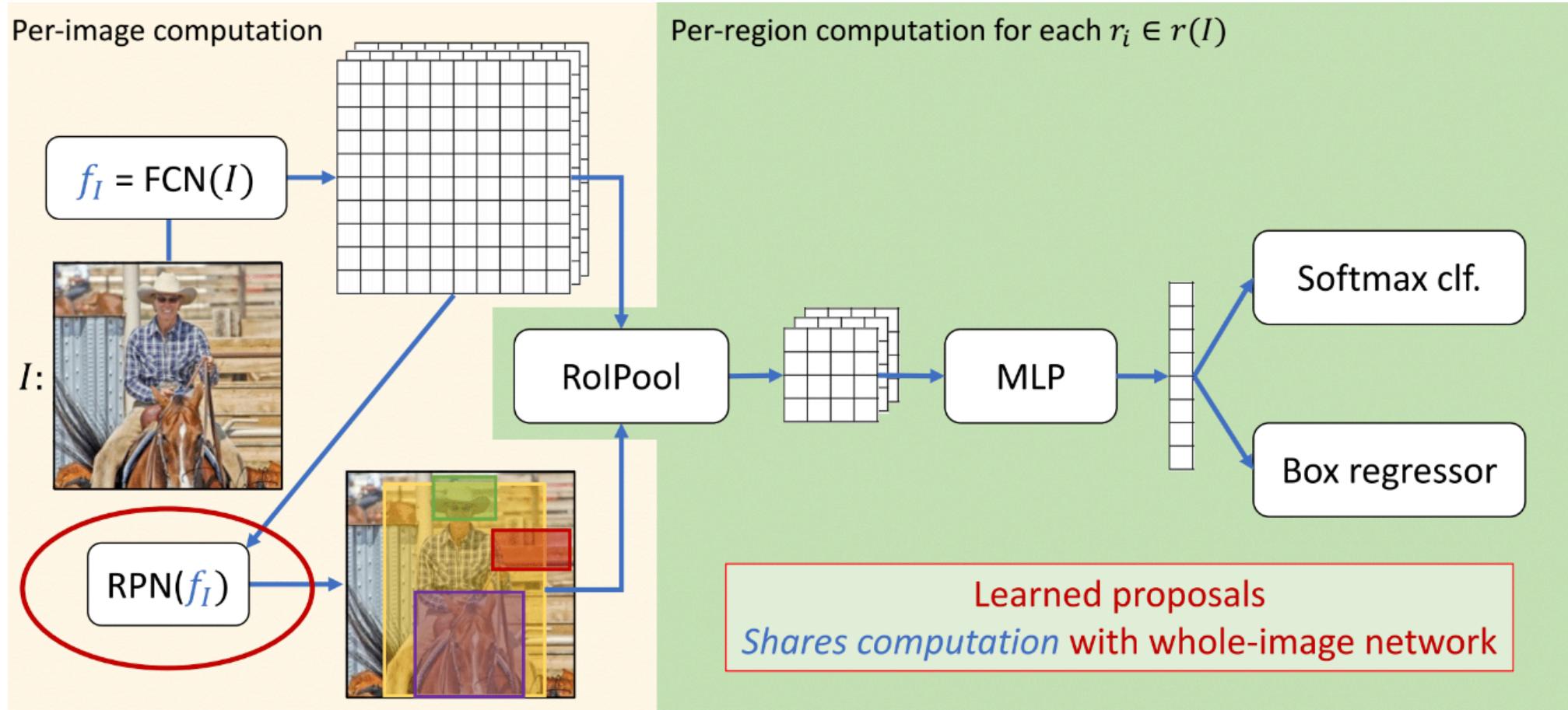
RoIPool on Each Proposal



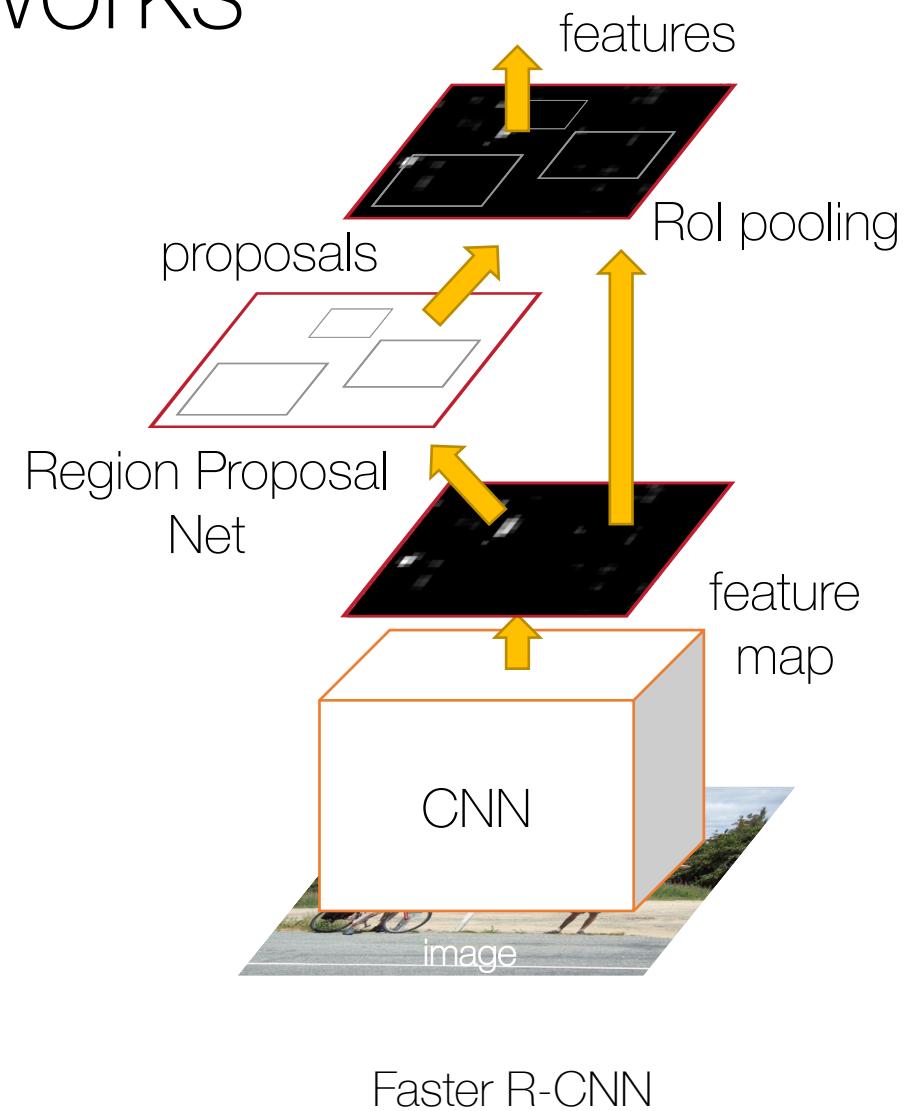
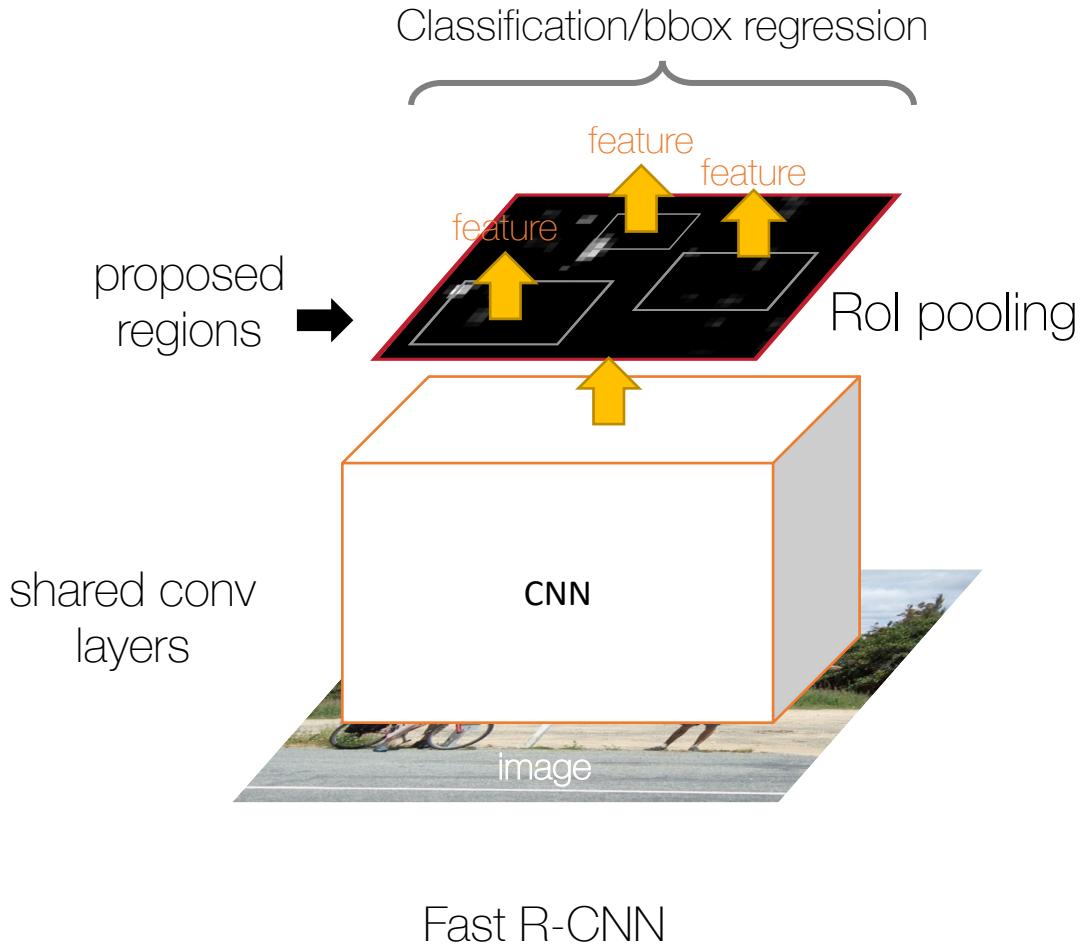
The Problem with Fast R-CNN



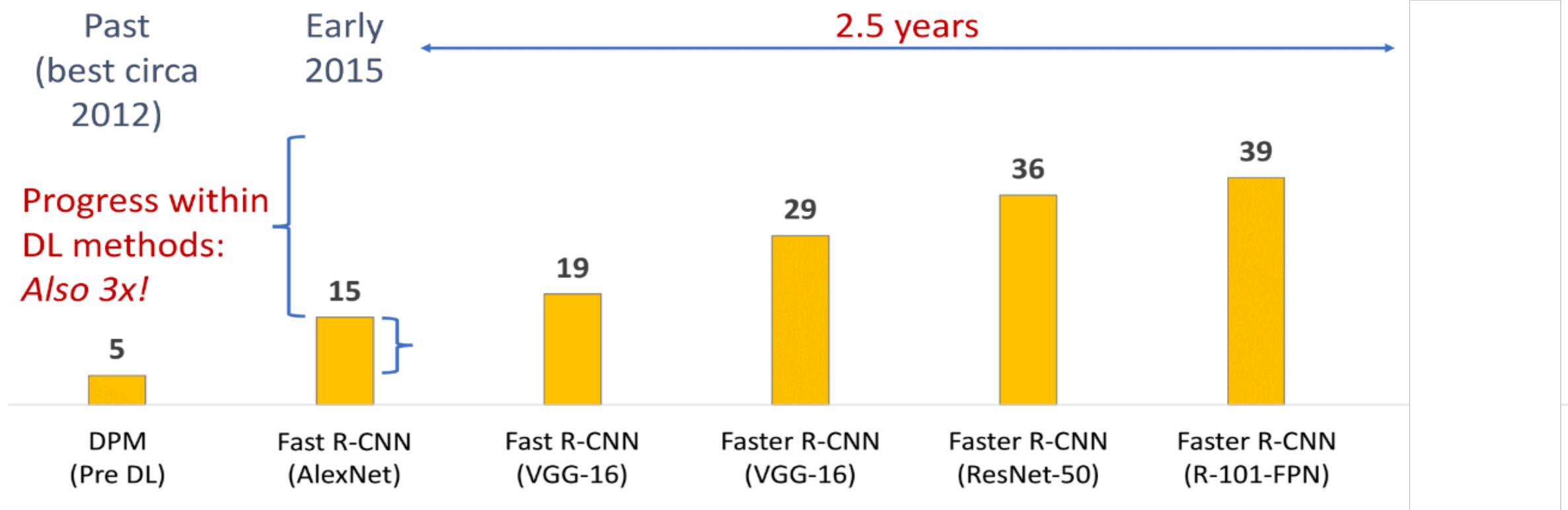
Object Detection Frameworks: Faster R-CNN



Object Detection Frameworks



Performance through Better Detector Design and Better Backbones



Instance Segmentation



Task

- Assume C object classes
- Segment all object instances in an image

True Positives

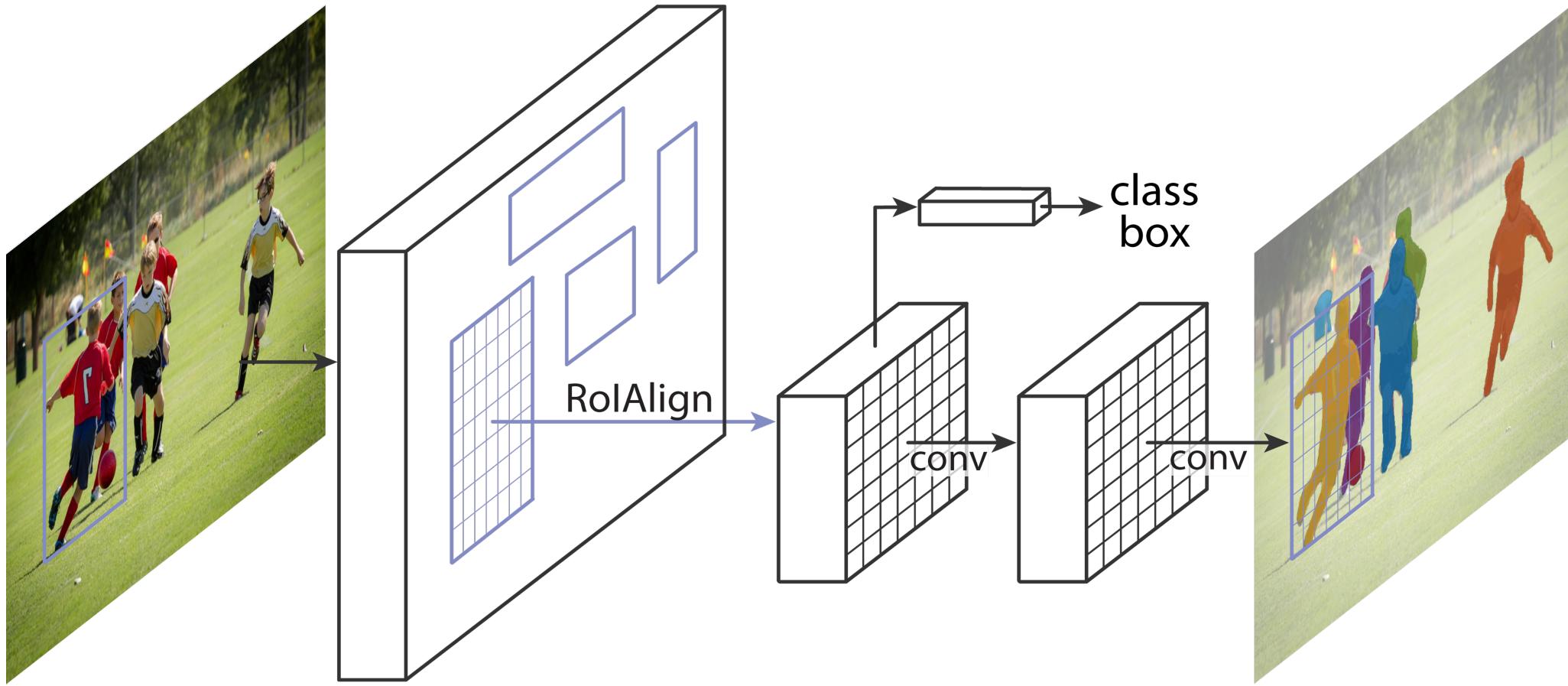
- A detection m is a true positive if $\text{maskov}(m, \text{gt}) > 0.5$ and $\text{class}(m) = \text{class}(\text{gt})$

False Positives

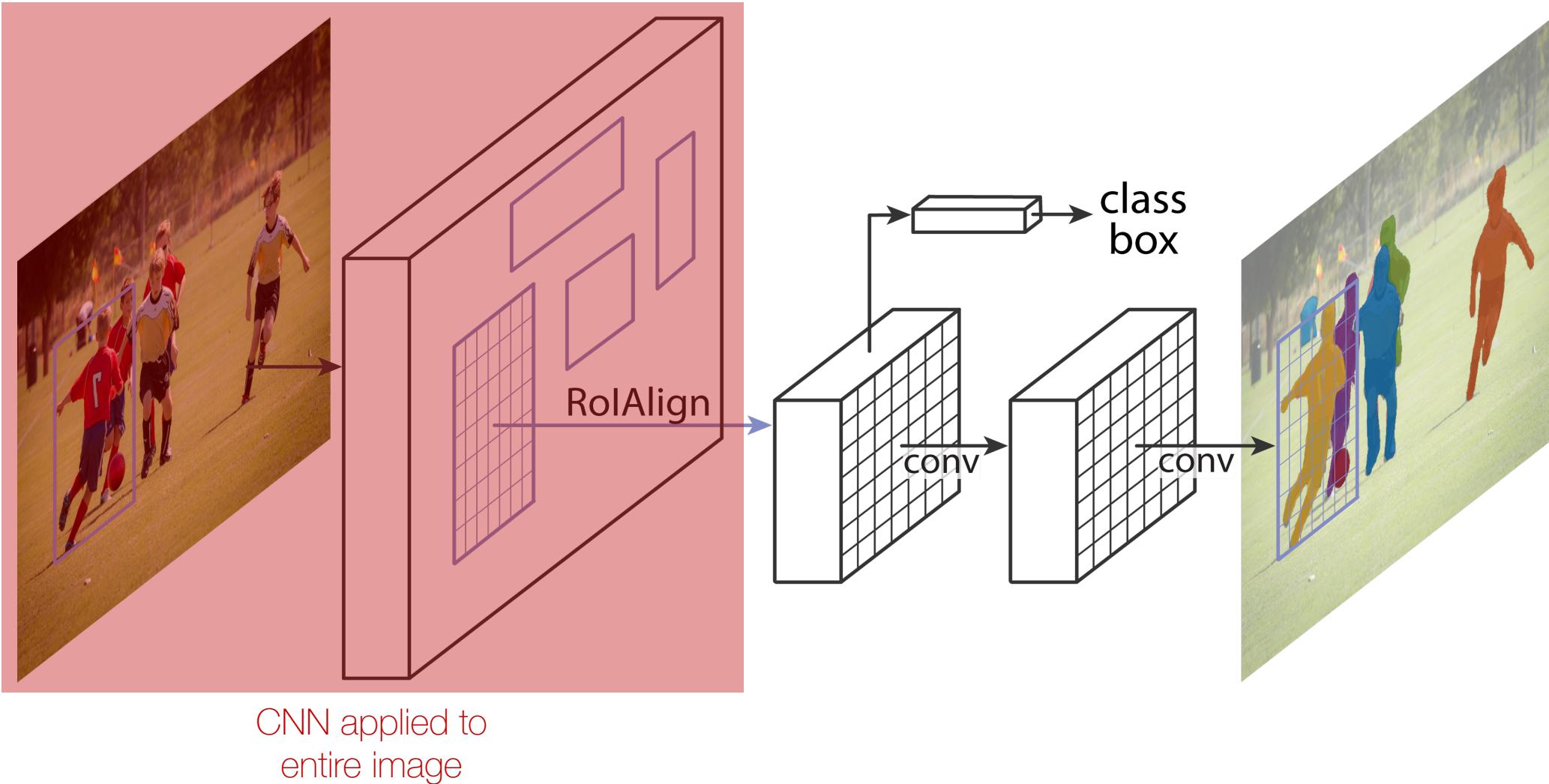
- Mislocalized and misclassified detections
- Duplicate detections are penalized

Mask R-CNN

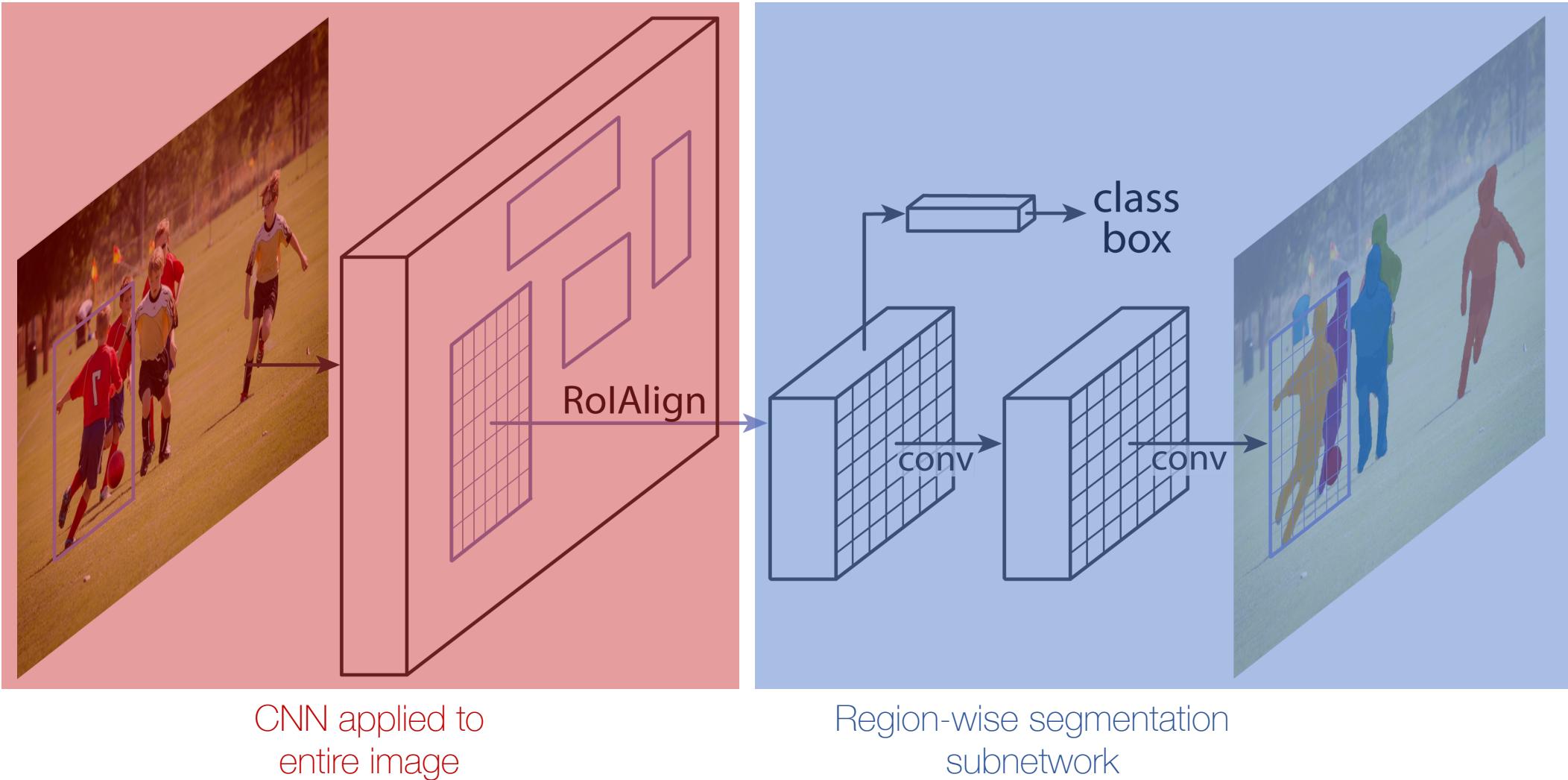
Mask R-CNN



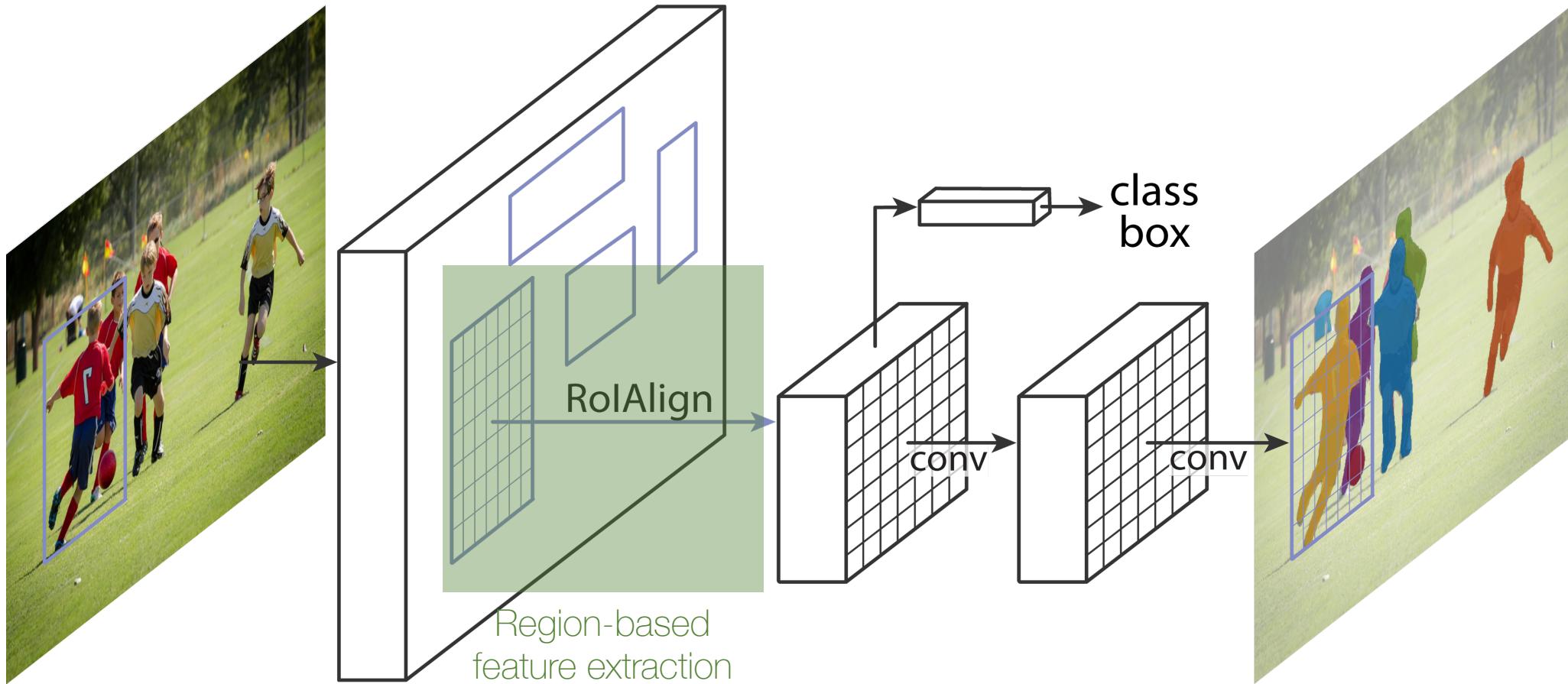
Mask R-CNN



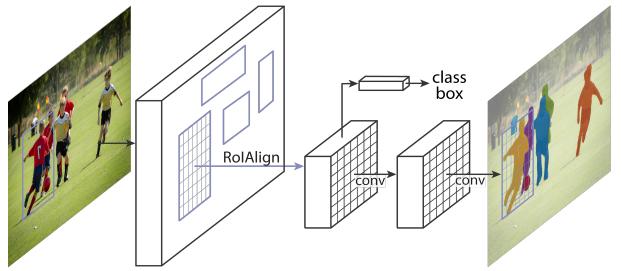
Mask R-CNN



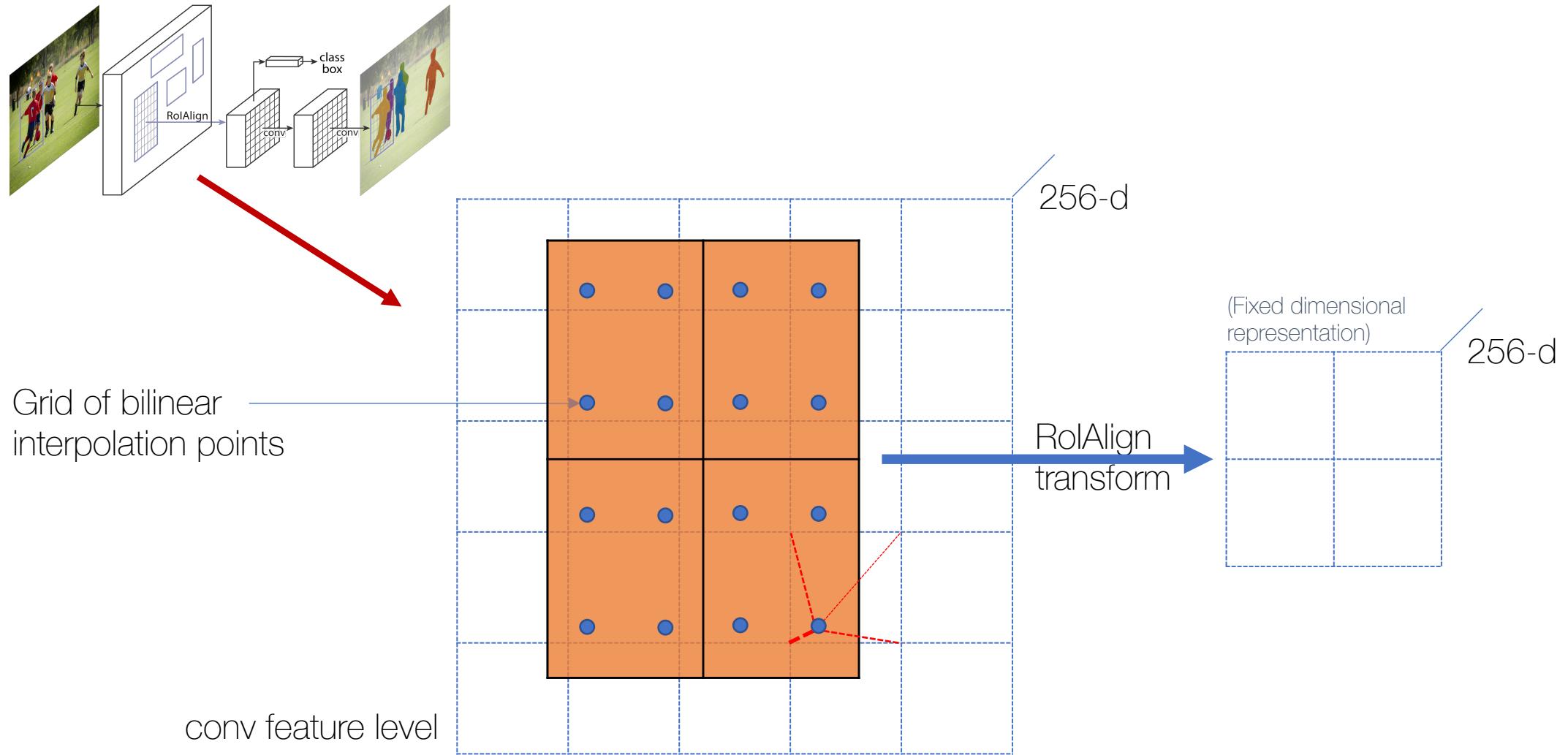
RoIAlign



RoIAlign

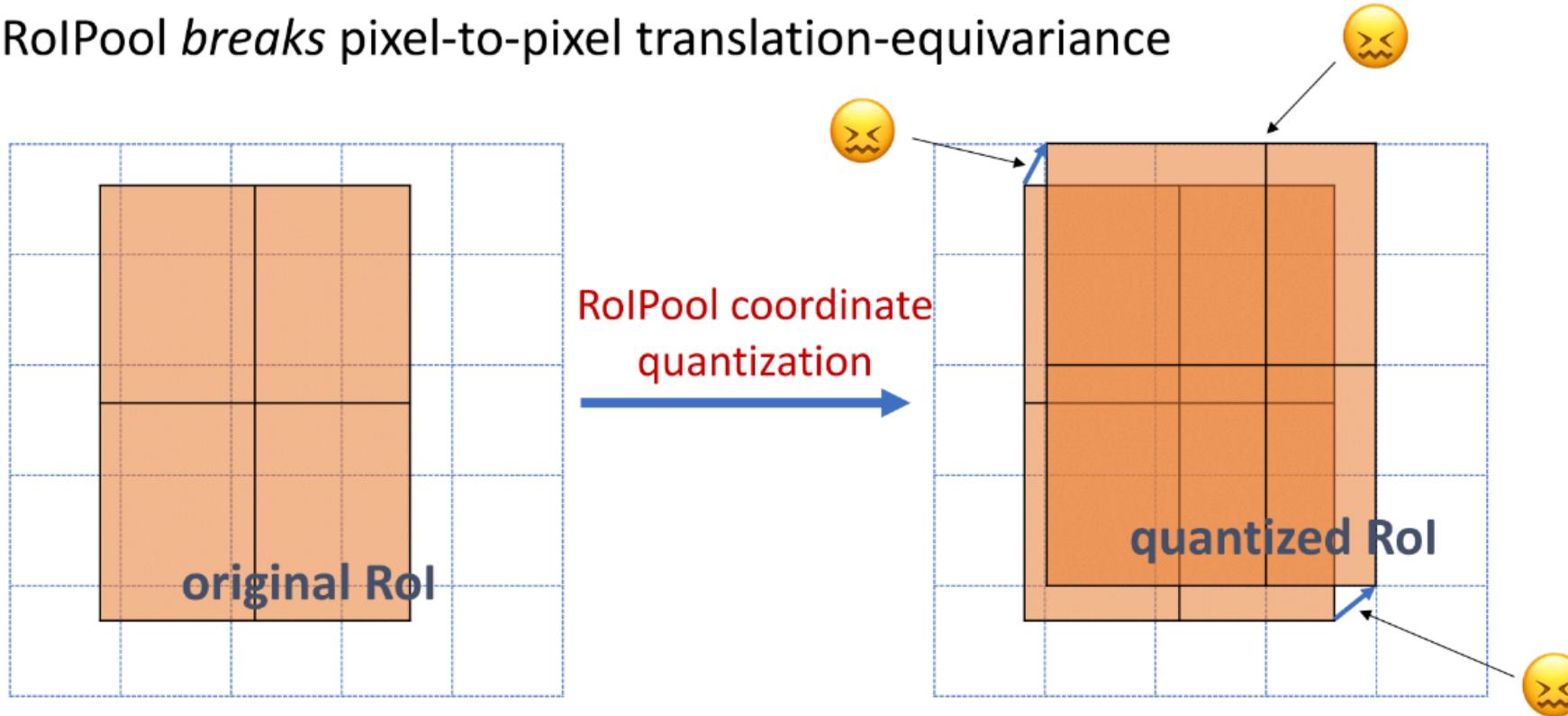


RoIAlign

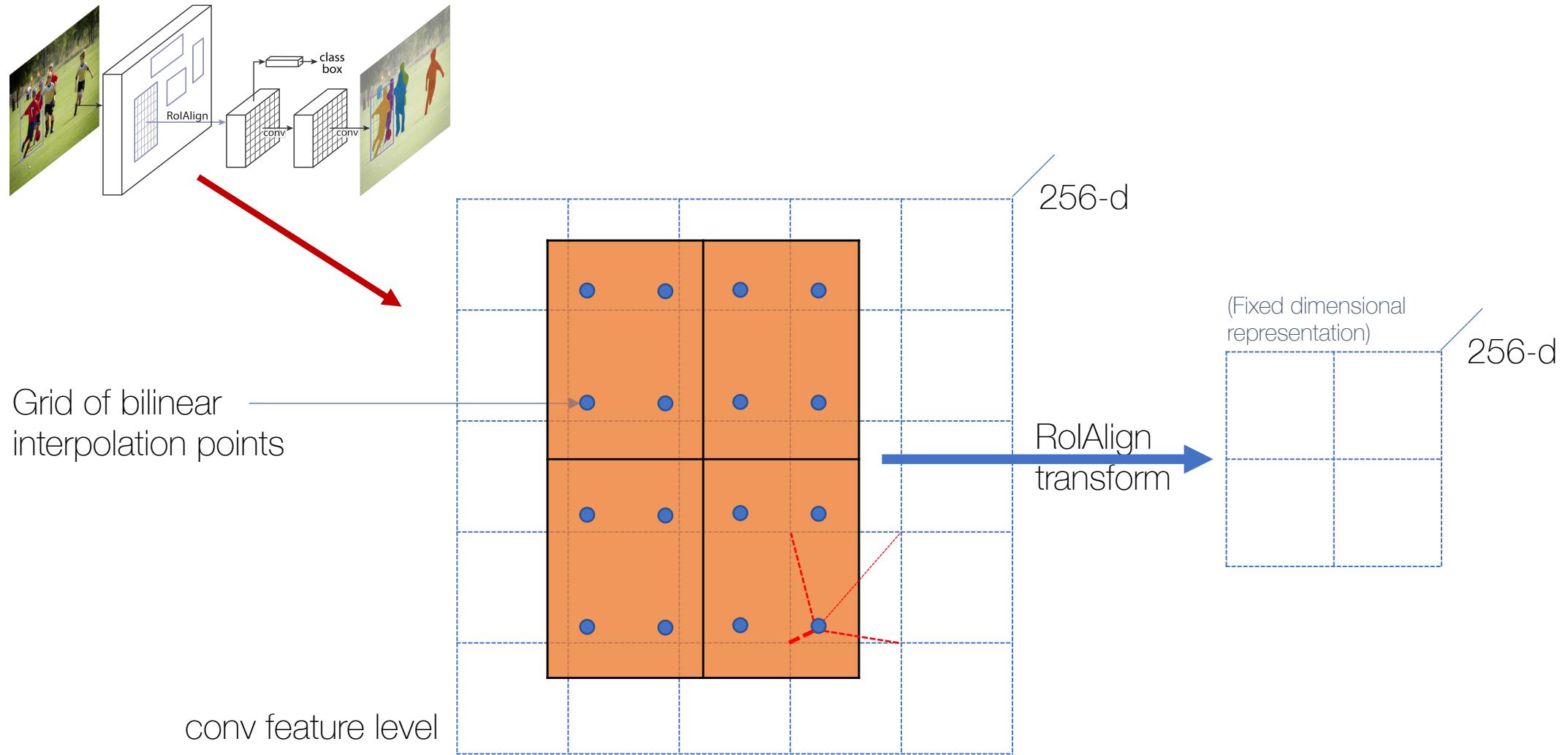


RoIPool vs RoIAlign

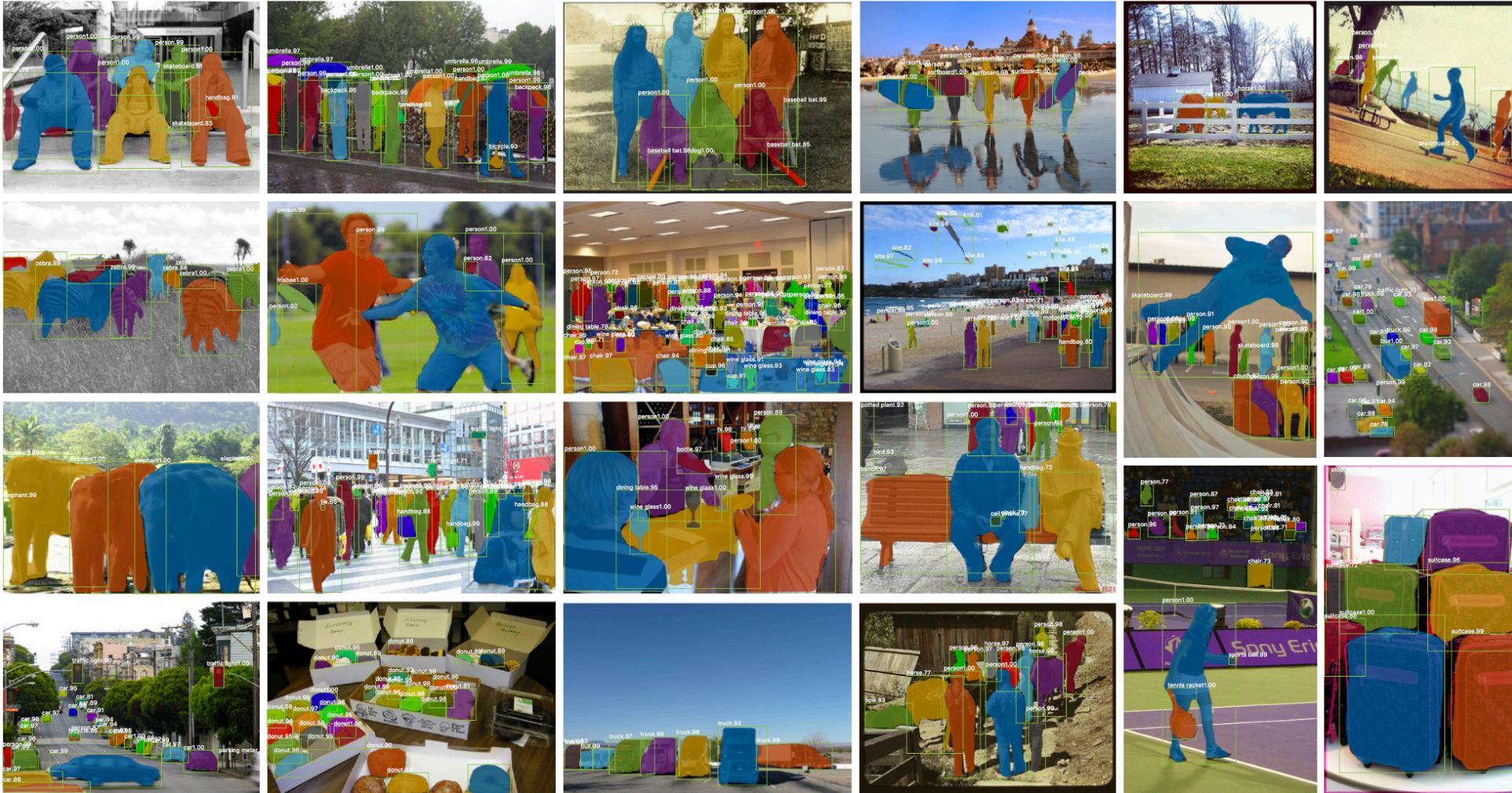
- RoIPool *breaks* pixel-to-pixel translation-equivariance



RoIAlign



Mask R-CNN for Instance Segmentation

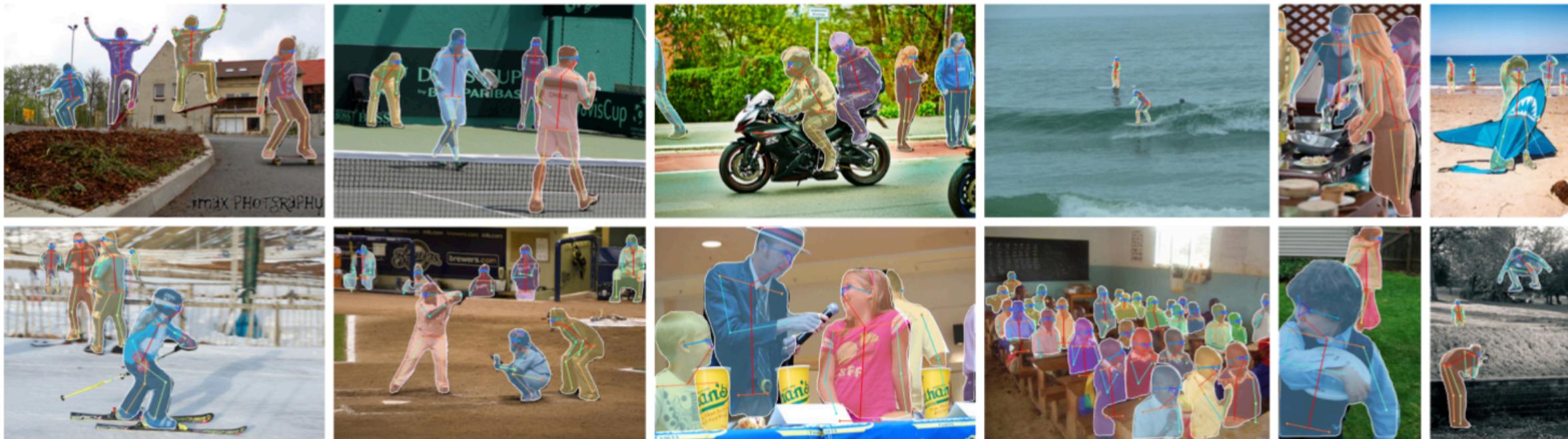


Mask R-CNN for Pose Prediction

Task

Given an image:

- Detect all human instances in the image
- Localize their landmarks, e.g. Right Shoulder, Nose etc.



Mask R-CNN for Pose Prediction

Loss

- Assume a set of K landmarks, e.g. Nose, Right Shoulder etc.
- For a detected instance R
 - corresponding logits J predicted by the Pose Head of Mask R-CNN of shape $H \times W \times K$
 - Ground truth locations $J_{gt} = \{(x_k, y_k), \text{ for } k = 0, \dots, K-1\}$

Content

- A Brief Intro
- Visual Recognition: Problem Definitions & Tasks
 - Semantic Segmentation
 - FCN
 - Object Detection
 - Fast(er) R-CNN
 - Instance Segmentation & Pose Prediction
 - Mask R-CNN
- ImageNet in 1H: Important Lessons
- Useful Tips

Making It Work

Solving a problem:

- Modeling: Transform the problem into a computational algorithm
- Architecture: Design the hierarchical model to capture the information necessary
- Invariance & Equivariance Properties
- Efficiency

Making It Work

Solving a problem:

- Modeling: Transform the problem into a computational algorithm
- Architecture: Design the hierarchical model to capture the information necessary
- Invariance & Equivariance Properties
- Efficiency
- **Hyperparameter Tuning**

Making It Work

Solving a problem:

- Modeling: Transform the problem into a computational algorithm
- Architecture: Design the hierarchical model to capture the information necessary
- Invariance & Equivariance Properties
- Efficiency
- **Hyperparameter Tuning: batch size, learning rate, epochs, decay, weight loss etc.**

Hyperparameters & SGD

- Batch Size: #gpus * #batch_per_gpu
- Learning Rate
- Training Schedule
- Weight Loss (important for multi-task learning)

Hyperparameters & SGD

- Batch Size: #gpus * #batch_per_gpu
- Learning Rate
- Training Schedule
- Weight Loss (important for multi-task learning)

These hyperparameters are correlated

ImageNet Training

Original Training Schedule [1]:

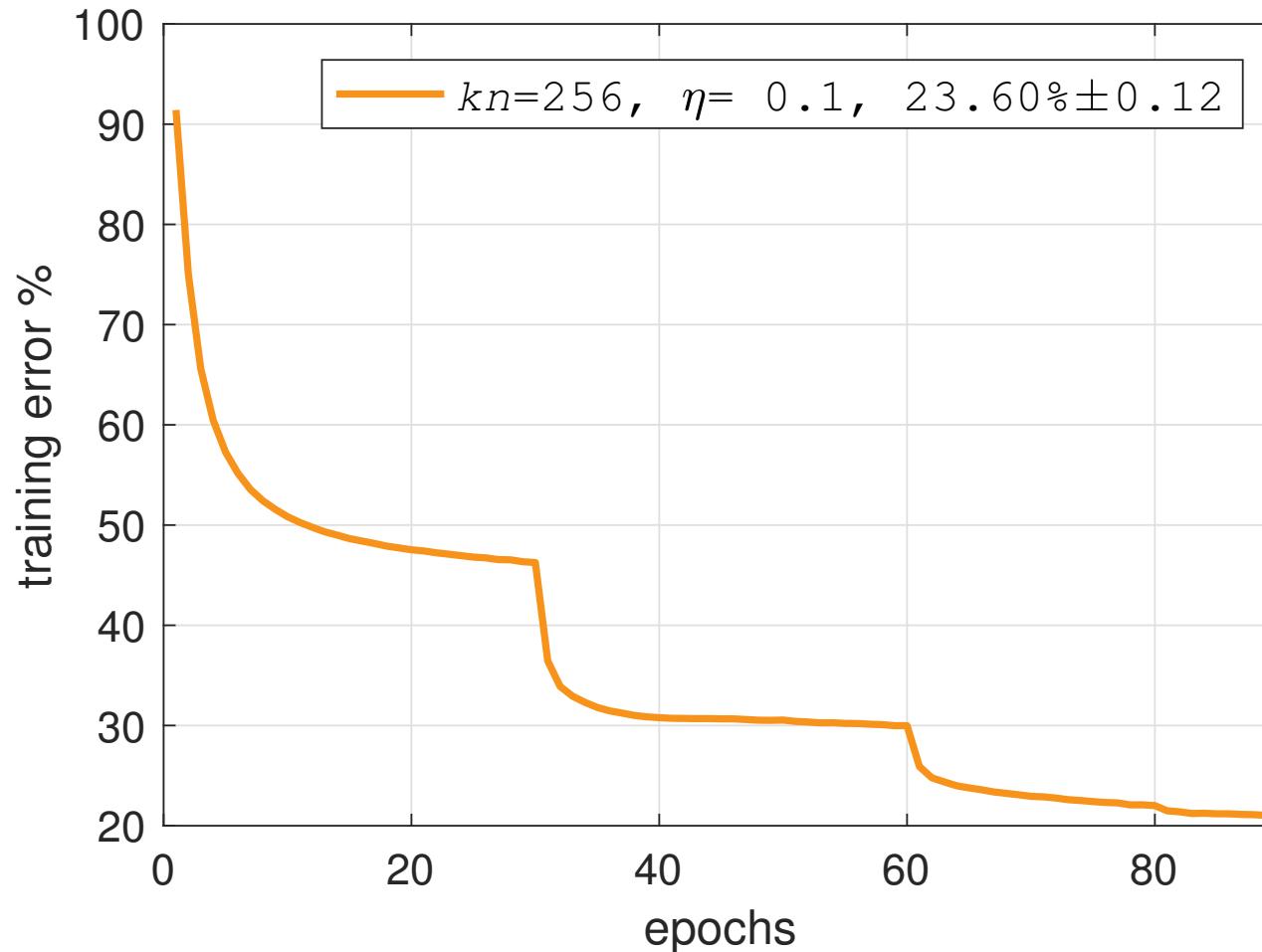
BASE_LR: 0.1

STEP_SIZES: [150150, 150150, 150150] # 30, 30, 30 epochs

LRS: [1, 0.1, 0.01, 0.001]

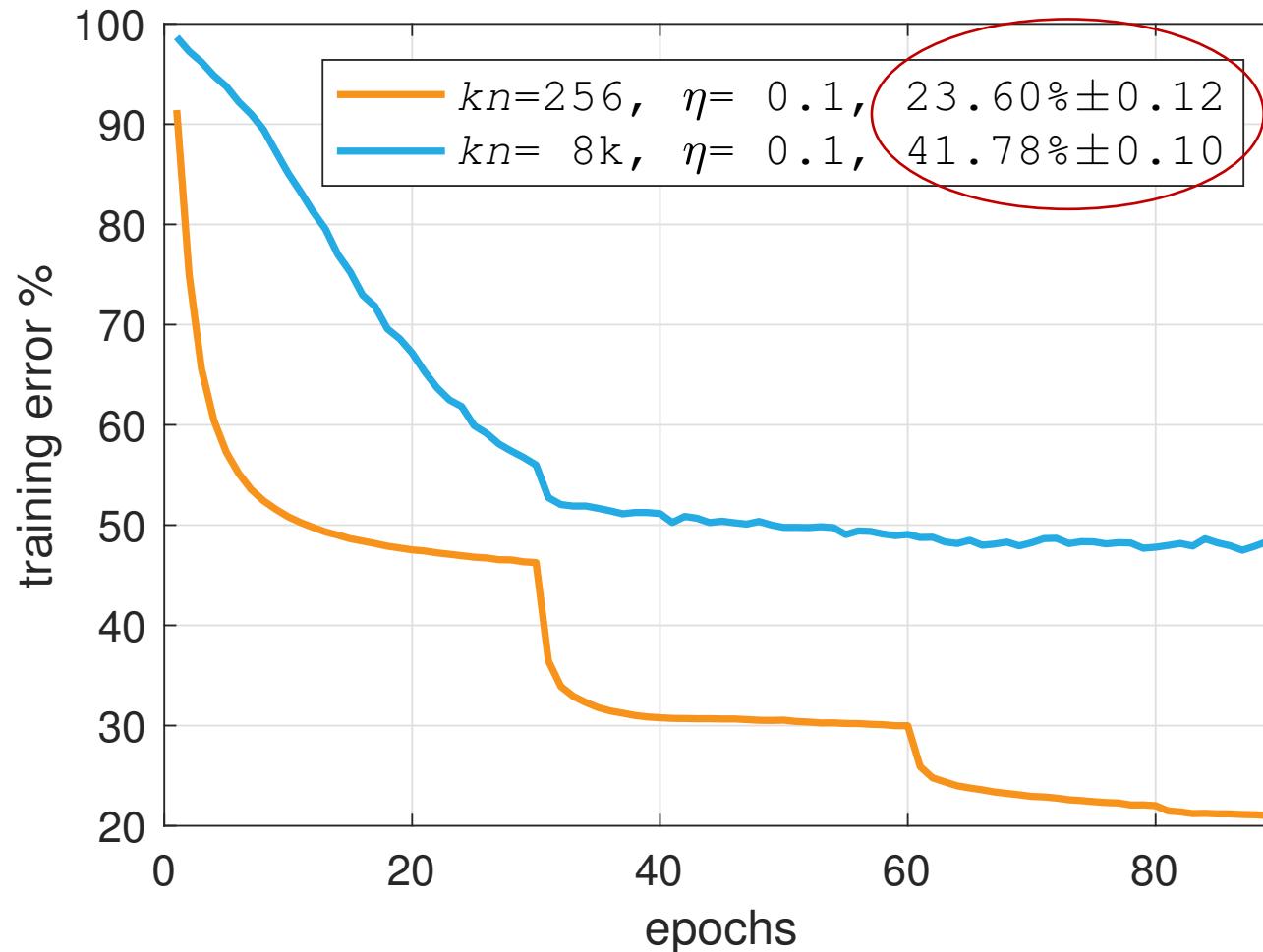
MAX_ITER: 500500 # total 100 epochs

ImageNet Training



- $k = \#gpus$
- $n = \text{per gpus batch size}$
- $\eta = \text{learning rate}$
- $256 = 8 \times 32$

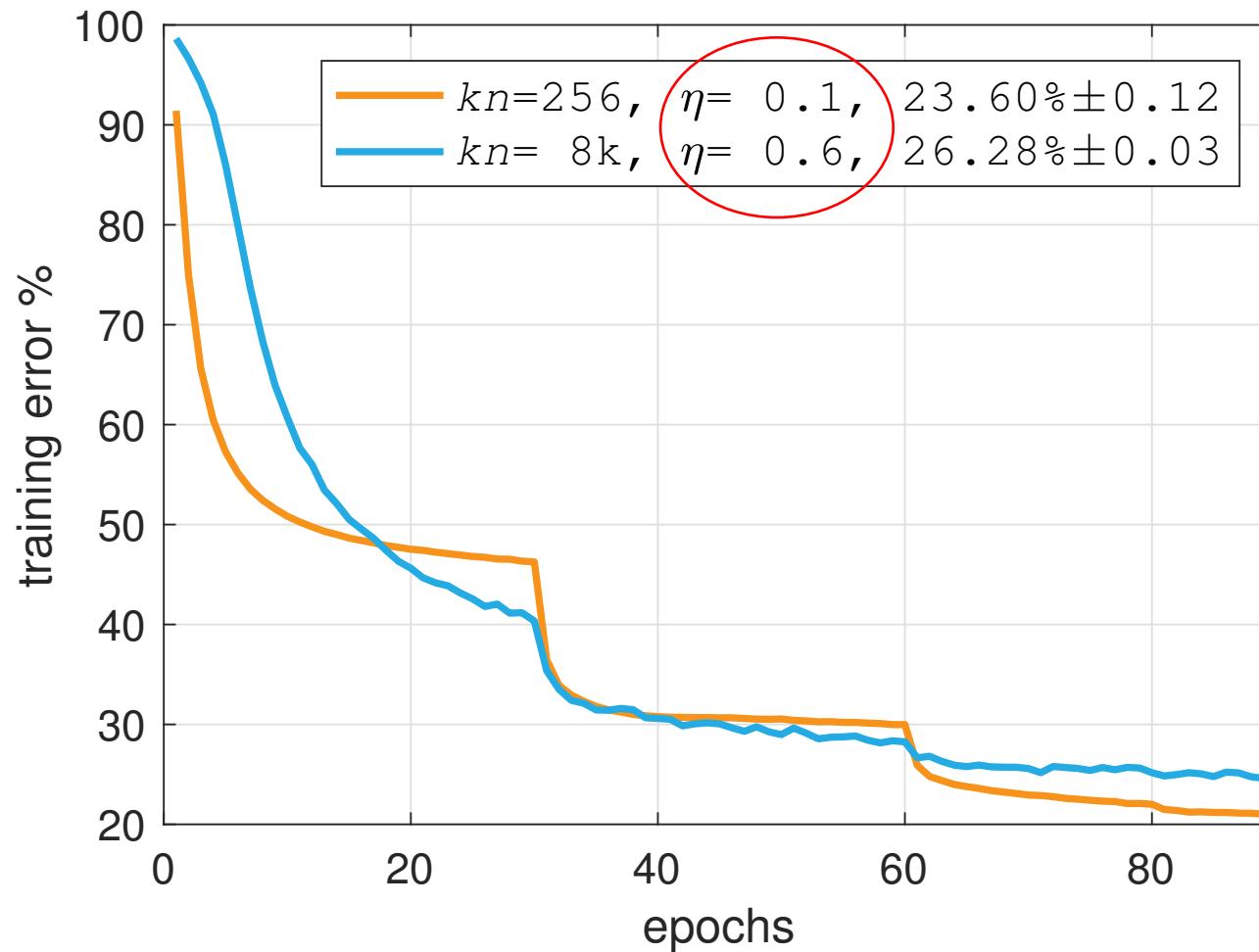
ImageNet in 1H [2]



$$8192 = 256 \times 32$$

↑
#gpus ↑
per gpu batch

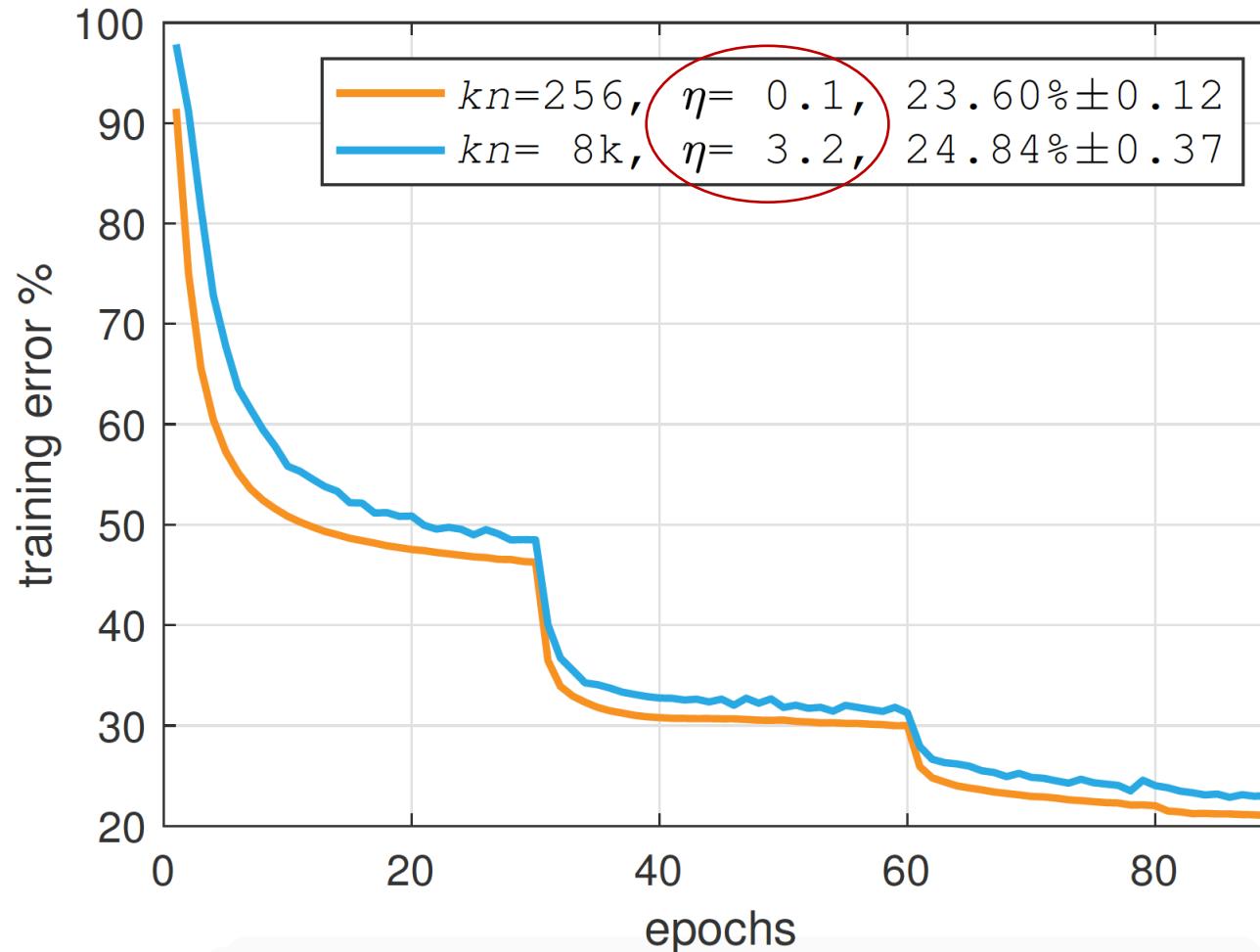
ImageNet in 1H [2]



Sqrt LR scaling

$$\eta \leftarrow \sqrt{\left(\frac{8192}{256}\right)} \cdot \eta$$

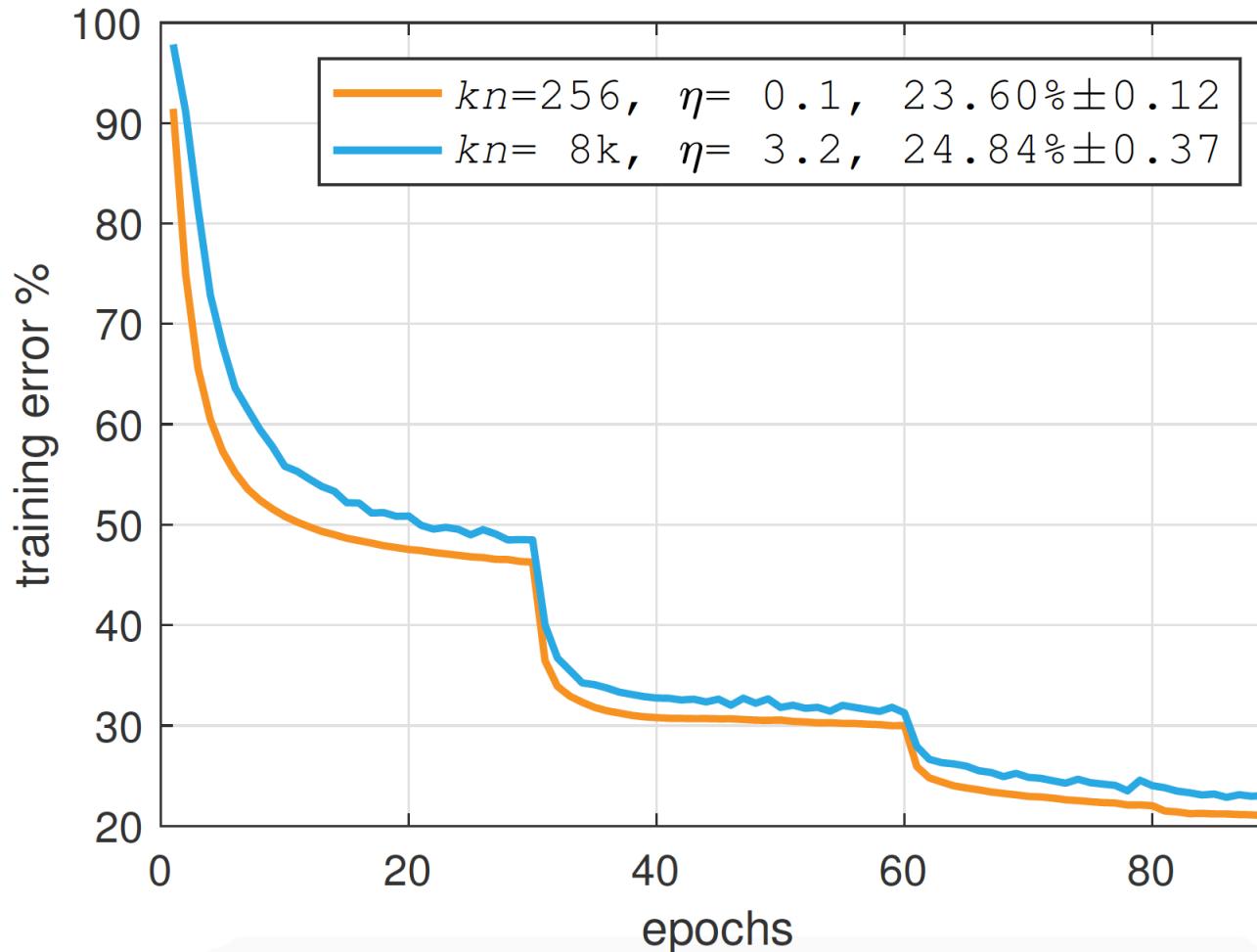
ImageNet in 1H [2]



Linear scaling

$$\eta \leftarrow \frac{8192}{256} \cdot \eta$$

ImageNet in 1H [2]

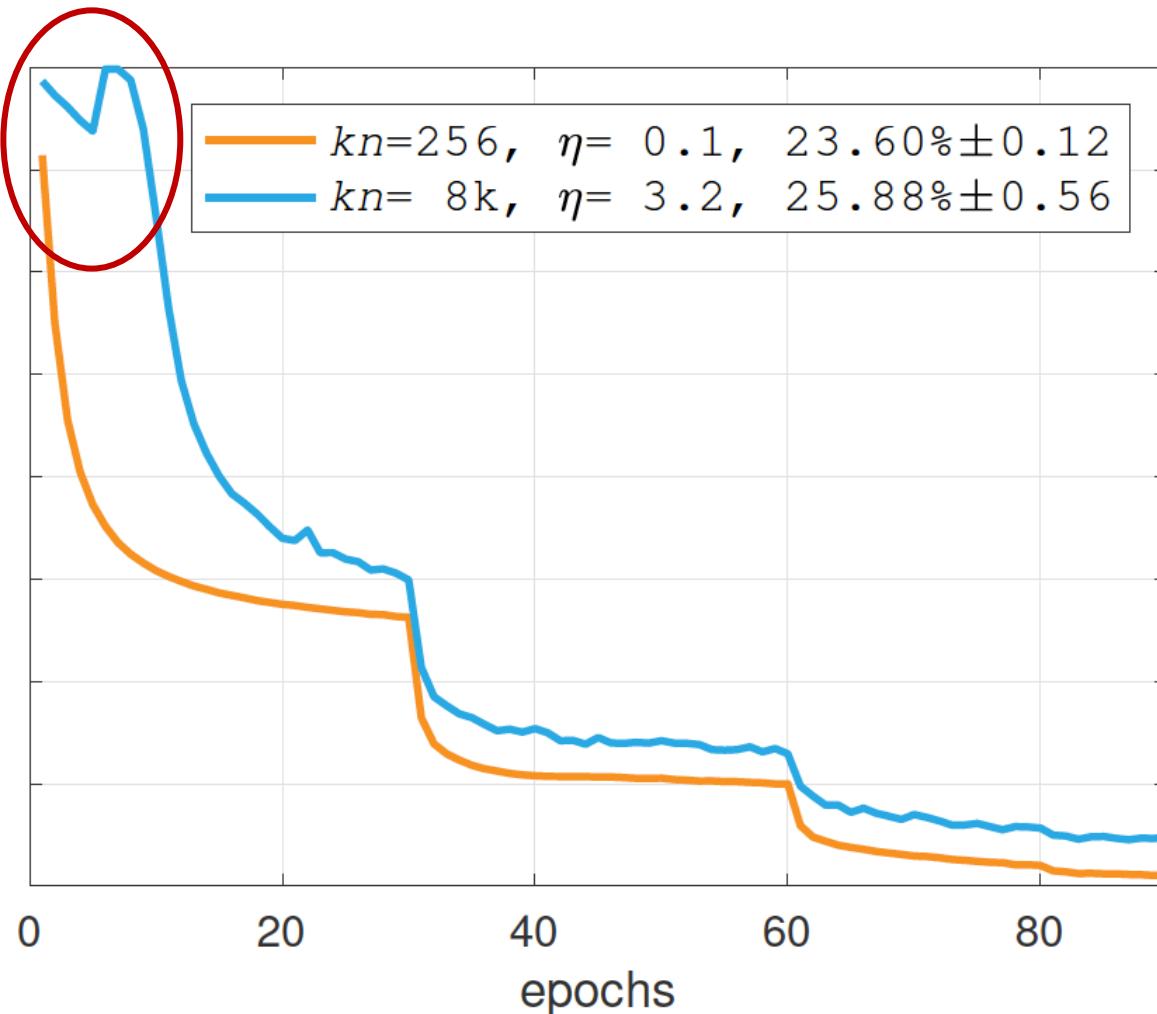


Linear scaling

$$\eta \leftarrow \frac{8192}{256} \cdot \eta$$

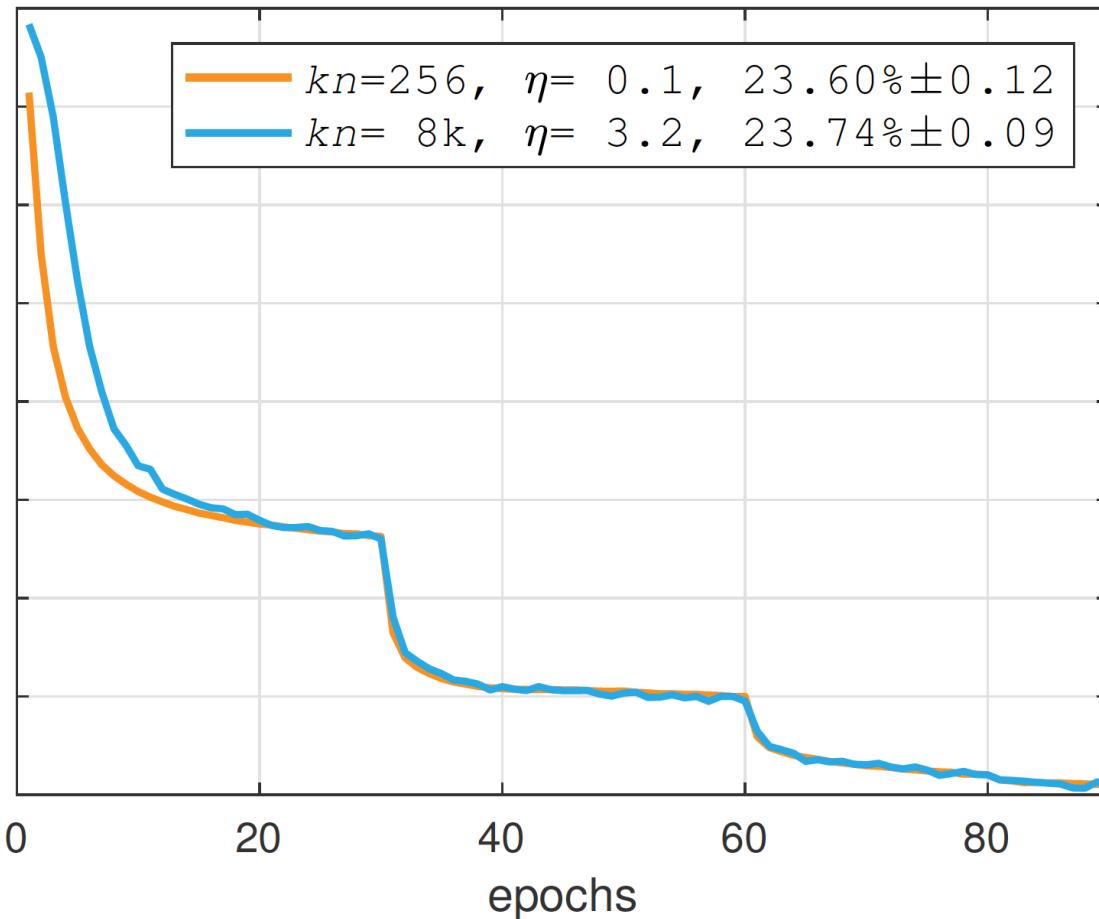
Optimization
difficulty

ImageNet in 1H [2]



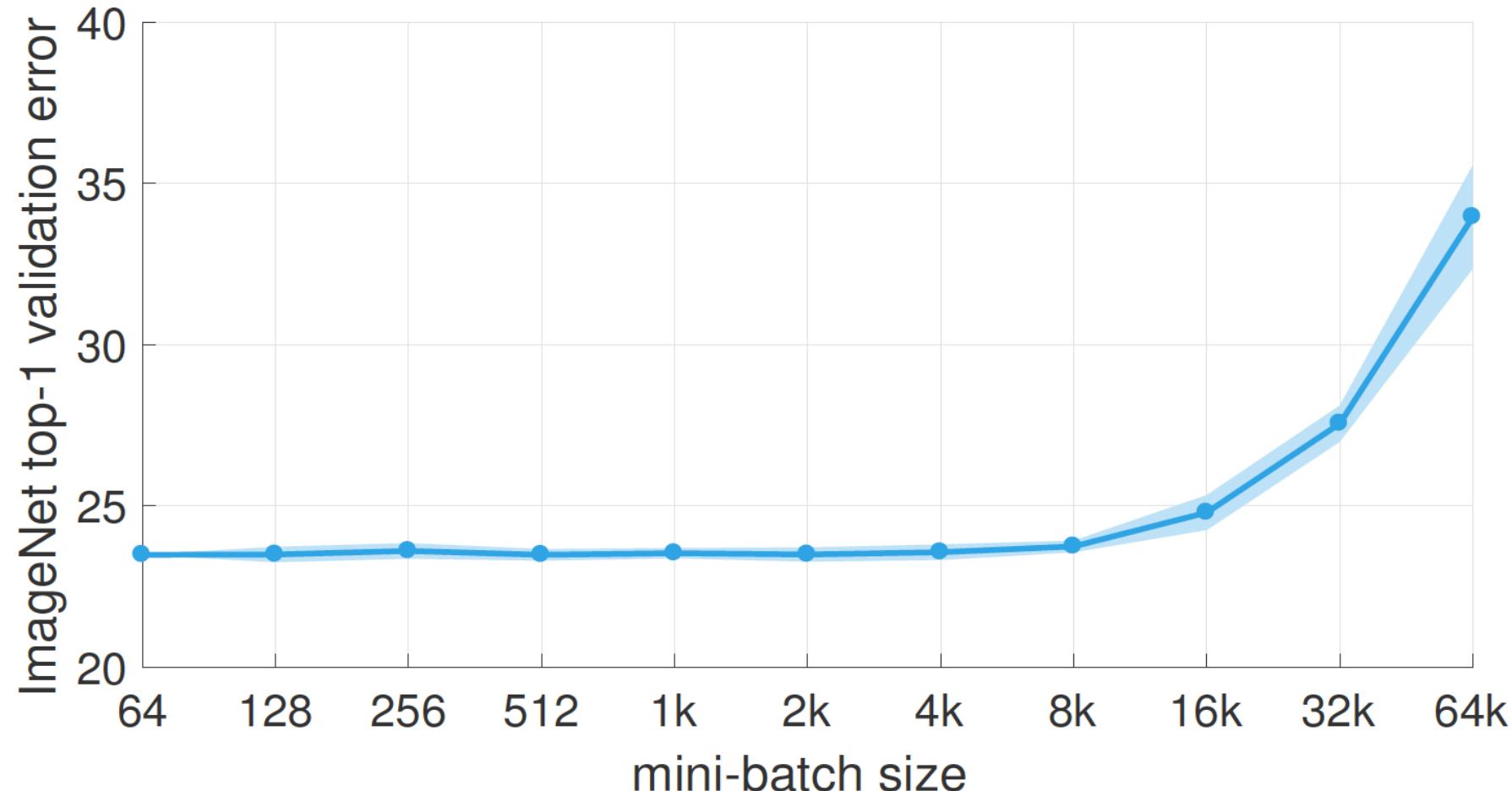
Rapid changes in network
in the beginning of training
→ use small LR for first few
epochs (**constant LR**
warmup)

ImageNet in 1H [2]



start from LR of $\eta(=0.1)$ and increase it by constant amount (0.003) at each iteration so that $\hat{\eta} = k\eta (= 3.2)$ after 5 epochs
(gradual LR warmup)

ImageNet in 1H [2]



[2] Goyal et al, Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour, 2017

Content

- A Brief Intro
- Visual Recognition: Problem Definitions & Tasks
 - Semantic Segmentation
 - FCN
 - Object Detection
 - Fast(er) R-CNN
 - Instance Segmentation & Pose Prediction
 - Mask R-CNN
- ImageNet in 1H: Important Lessons
- Useful Tips

Common Mistakes

- Remember that if you reduce the effective batch size you need to reduce the learning rate

Common Mistakes

- Remember that if you reduce the effective batch size you need to reduce the learning rate
- When using #gpus >1 be aware of how the gradients are combined (added vs. averaged) → tune learning rate and loss weight accordingly

```
kps_prob, loss_kps = model.net.SoftmaxWithLoss(  
    ['kps_score_reshaped', 'keypoint_locations_int32', 'keypoint_weights'],  
    ['kps_prob', 'loss_kps'],  
    scale=cfg.KRCNN.LOSS_WEIGHT / cfg.NUM_GPUS,  
    spatial=0  
)
```

Common Mistakes

- Remember that if you reduce the effective batch size you need to reduce the learning rate
- When using `#gpus > 1` be aware of how the gradients are combined (added vs. averaged) → tune learning rate and loss weight accordingly
- Count training duration in epochs not in training iterations: two models are comparable when they have been trained for the same number of epochs (this might not be the same number of iterations)

How to do research

- Read old & new papers
 - Identify what you are passionate about
 - Start from the basics and build up (instead of the opposite!)
 - Ablate all the parts of your approach
 - Visualize your input, output, intermediate representations
 - Do not use deep learning as a black box
 - Understand the **forward** flow of computation
 - Understand the **backward** flow of computation
- (see: [Andrej's blog post](#))